



# *University of* **HUDDERSFIELD**

## **University of Huddersfield Repository**

Lee, Hyunkook

Capturing and Rendering 360° VR Audio Using Cardioid Microphones

### **Original Citation**

Lee, Hyunkook (2016) Capturing and Rendering 360° VR Audio Using Cardioid Microphones. In: AES Conference on Audio for Augmented and Virtual Reality, 30 Sep - 1 Oct 2016, Los Angeles, USA.

This version is available at <http://eprints.hud.ac.uk/id/eprint/29582/>

The University Repository is a digital collection of the research output of the University, available on Open Access. Copyright and Moral Rights for the items on this site are retained by the individual author and/or other copyright owners. Users may access full items free of charge; copies of full text items generally can be reproduced, displayed or performed and given to third parties in any format or medium for personal research or study, educational or not-for-profit purposes without prior permission or charge, provided:

- The authors, title and full bibliographic details is credited in any copy;
- A hyperlink and/or URL is included for the original metadata page; and
- The content is not changed in any way.

For more information, including our policy and submission procedure, please contact the Repository Team at: [E.mailbox@hud.ac.uk](mailto:E.mailbox@hud.ac.uk).

<http://eprints.hud.ac.uk/>

# **Capturing and Rendering 360° VR Audio using Cardioid Microphones**

Hyunkook Lee

[h.lee@hud.ac.uk](mailto:h.lee@hud.ac.uk)

Applied Psychoacoustics Lab (APL)  
University of Huddersfield, UK

---

- Near-coincident mic arrays
  - ORTF, NOS, etc.
  - Arguably, preferred to pure coincident or pure spaced techniques by most professional recording engineers.
  - Rely on the trade-off between Time and Level differences.
  - Best of both worlds (Localisability & Spaciousness).
- Cardioid microphones
  - Most popular.
  - Most widely available.
- Record for VR using favourite cardioid mics arranged in a near-coincident fashion?

- Research background
- Localisation test in loudspeaker reproduction
- Localisation test in binaural reproduction
- Discussion
- Summary

# Research Background



# Existing methods for VR audio capture

- First Order Ambisonics (FOA)



Pros	Cons
<ul style="list-style-type: none"><li>• Very good “localisability” due to the coincident nature (But not necessarily good localisation “accuracy”).</li><li>• Virtual microphones from flexible decoding.</li><li>• Compact.</li></ul>	<ul style="list-style-type: none"><li>• High interchannel correlation.</li><li>• Lack of spaciousness.</li><li>• Comb-filtering and rapid change in image position even with a small head movement.</li></ul>

# Existing methods for VR audio capture

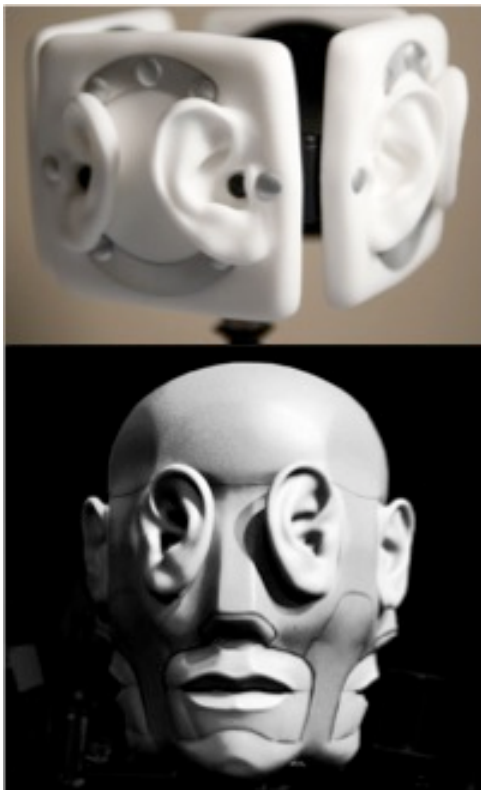
- Higher Order Ambisonics (HOA)



Pros	Cons
<ul style="list-style-type: none"><li>• Higher spatial resolution.</li><li>• More accurate localisation.</li></ul>	<ul style="list-style-type: none"><li>• Requires a large number of channels for a proper decoding. <math>N = (M + 1)^2</math></li><li>• Very expensive.</li><li>• Tonal quality.</li><li>• Spaciousness?</li></ul>

# Existing methods for VR audio capture

- Quad Binaural

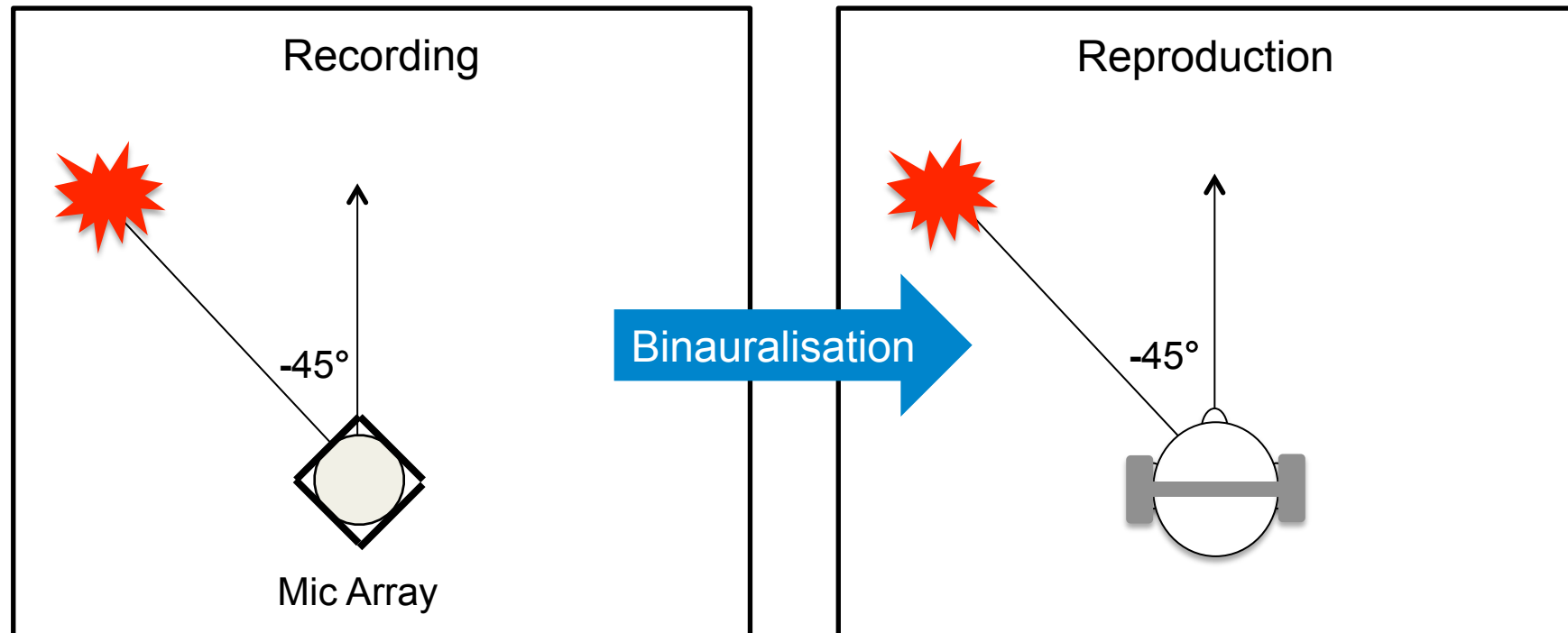


Pros	Cons
<ul style="list-style-type: none"><li>• Direct pinnae filtering.</li><li>• No need for extra binaural synthesis.</li></ul>	<ul style="list-style-type: none"><li>• Inaccurate localisation and comb-filtering due to crossfading between ear signals.</li><li>• Not possible to use personal HRTFs.</li><li>• Only for horizontal head rotation.</li><li>• Expensive.</li></ul>



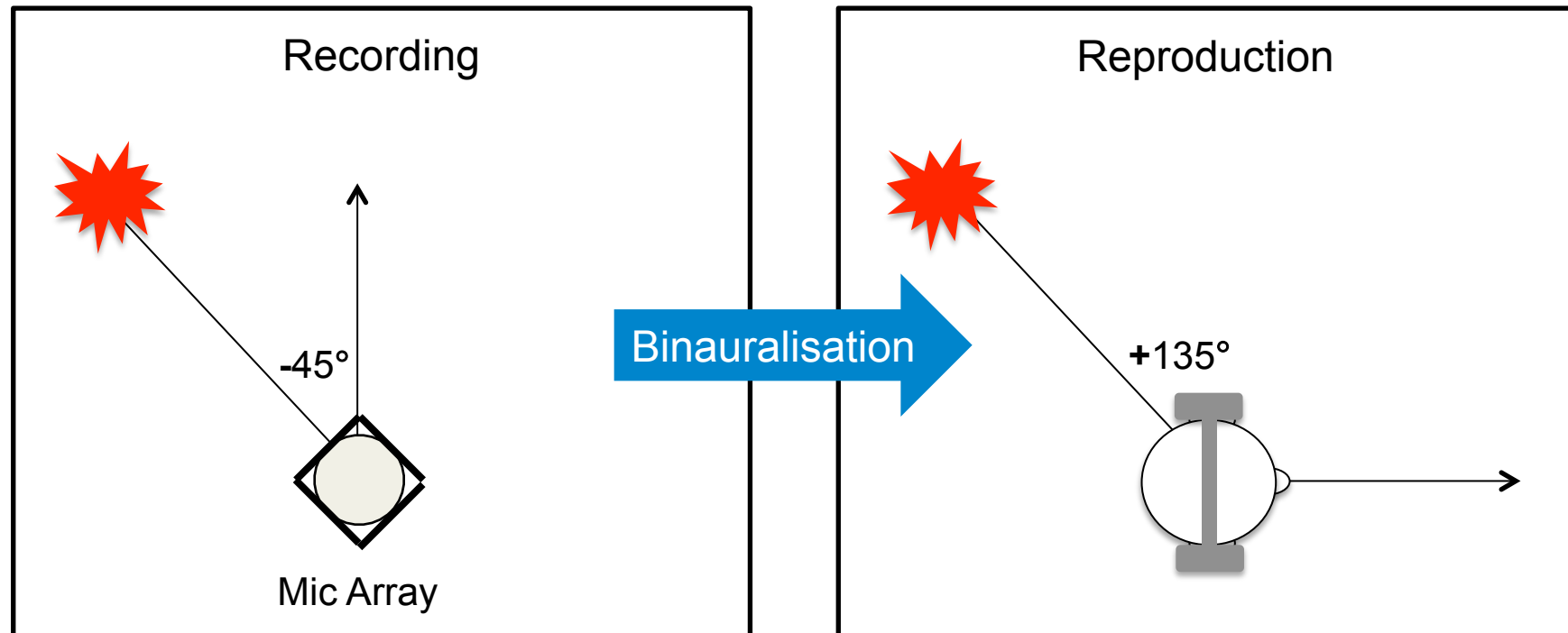
# Psychoacoustic considerations for VR

- In VR, it is important to match the actual and perceived source positions.



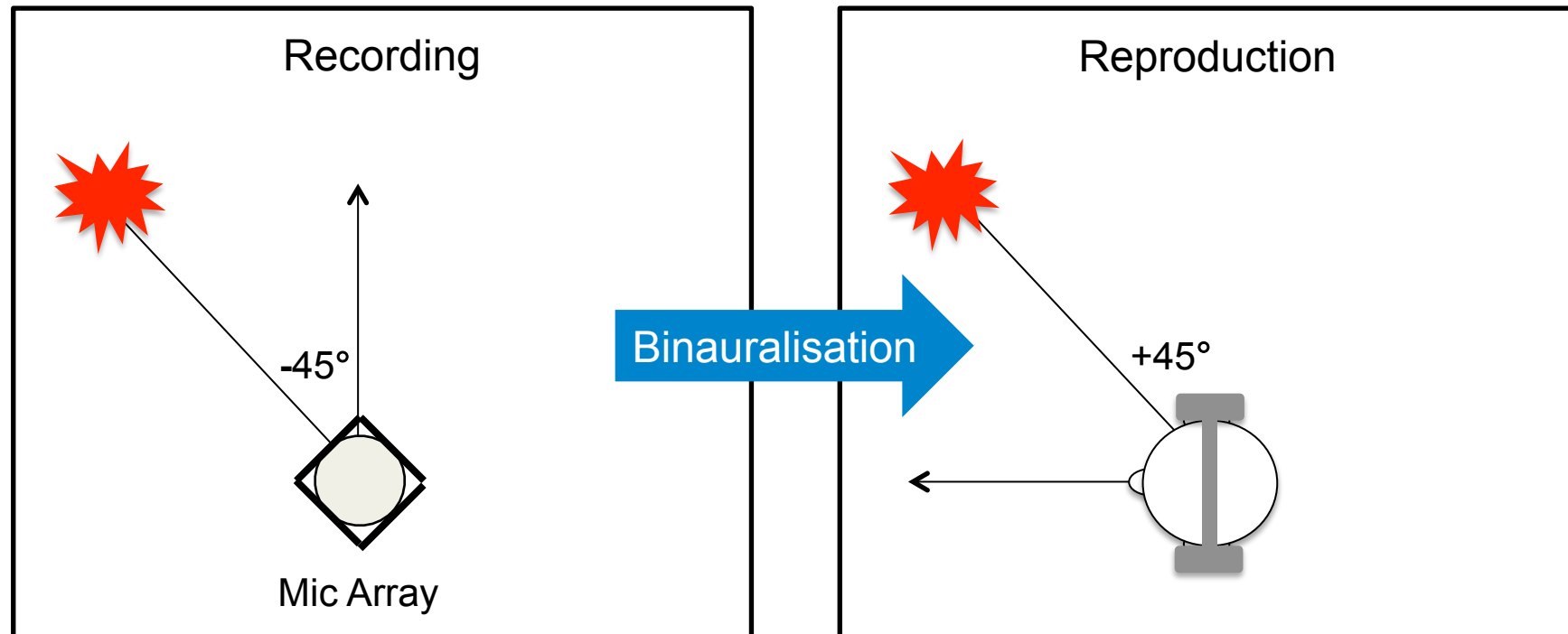
# Psychoacoustic considerations for VR

- The perceived source position should stay the same as the head rotates.



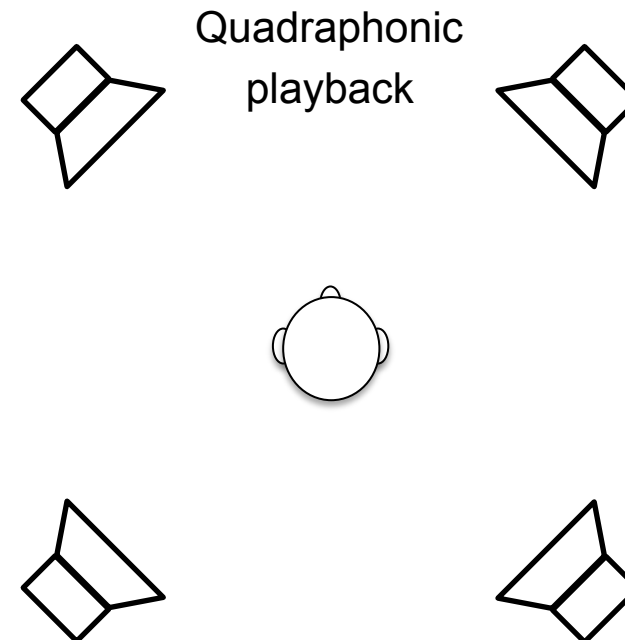
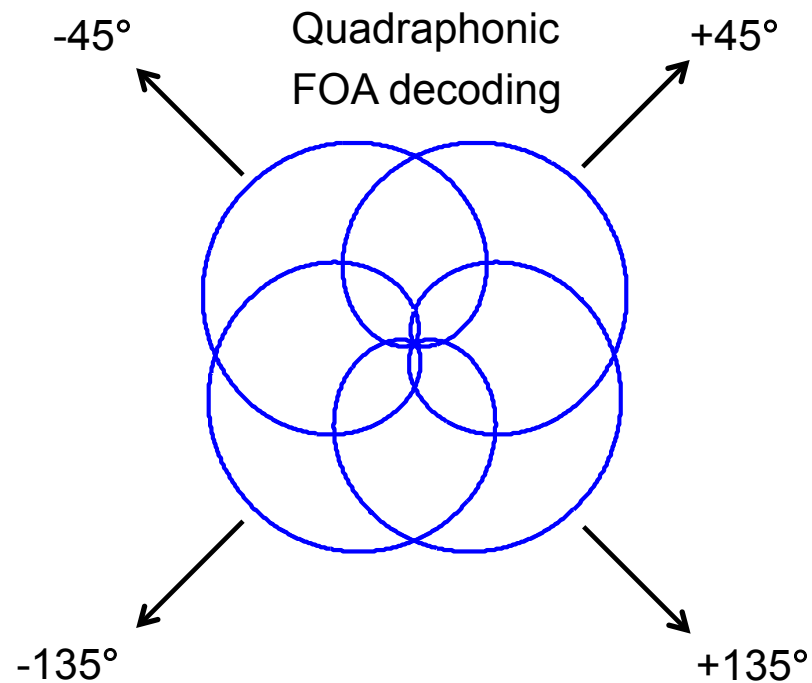
# Psychoacoustic considerations for VR

- The perceived source position should stay the same as the head rotates.



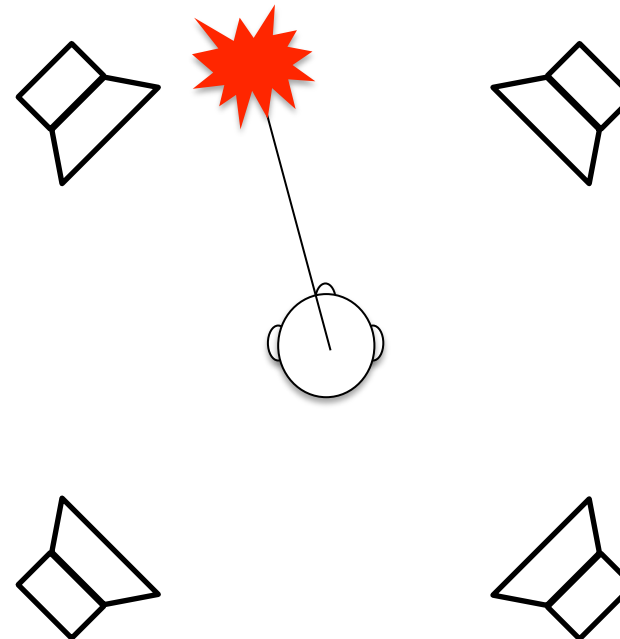
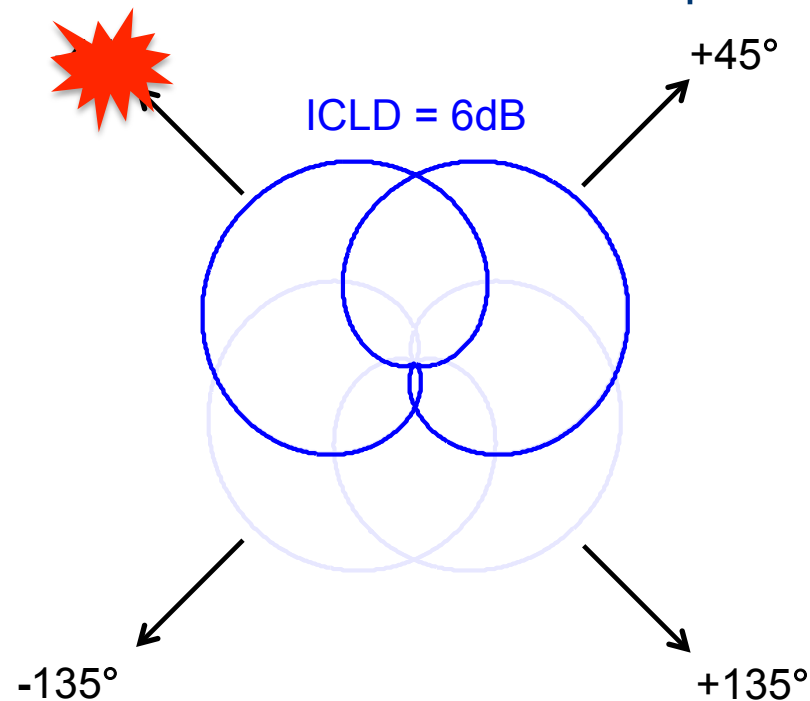
# Psychoacoustic considerations for VR

- Limitation of FOA
  - Quadraphonic Cardioid decoding.



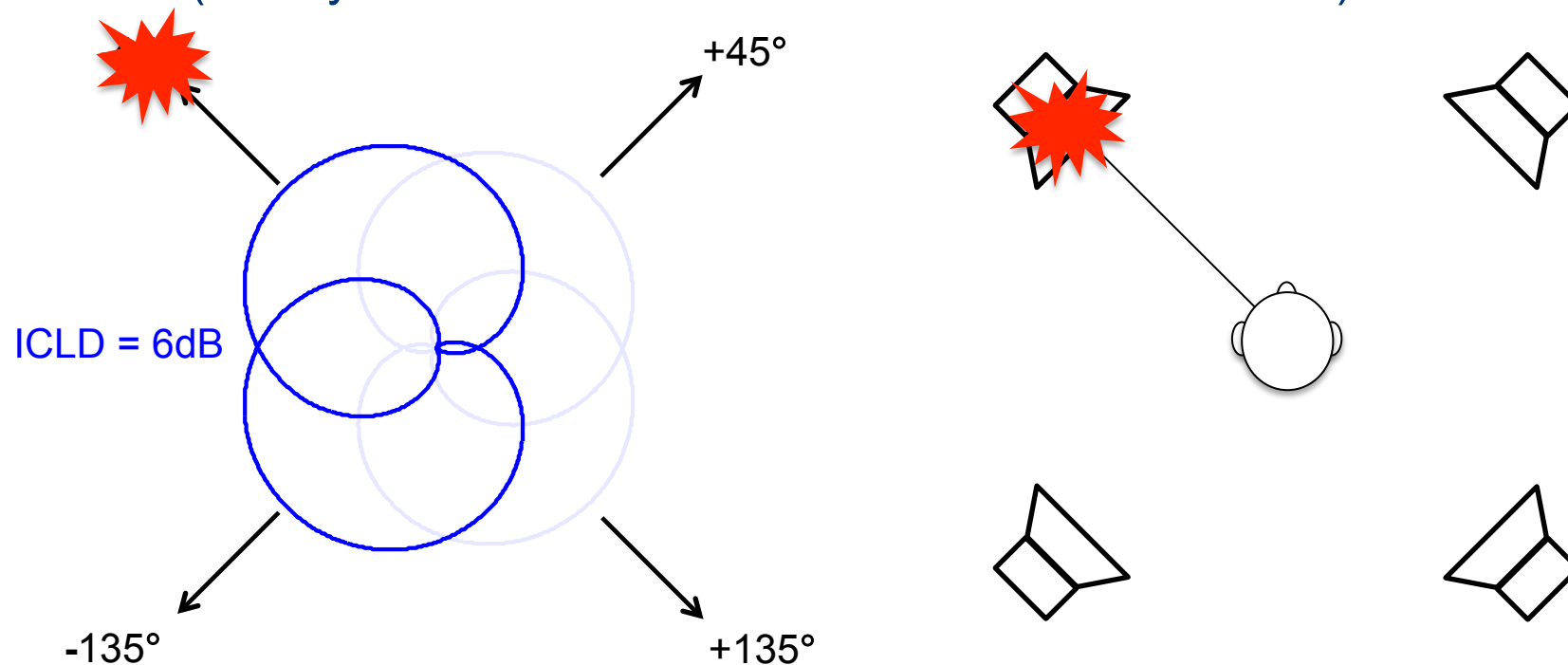
# Psychoacoustic considerations for VR

- Limitation of FOA
  - Only **6dB ICLD** (interchannel level difference) for the front pair for a source at  $45^\circ$ .
  - Not sufficient for a full phantom image shift to  $45^\circ$ .



# Psychoacoustic considerations for VR

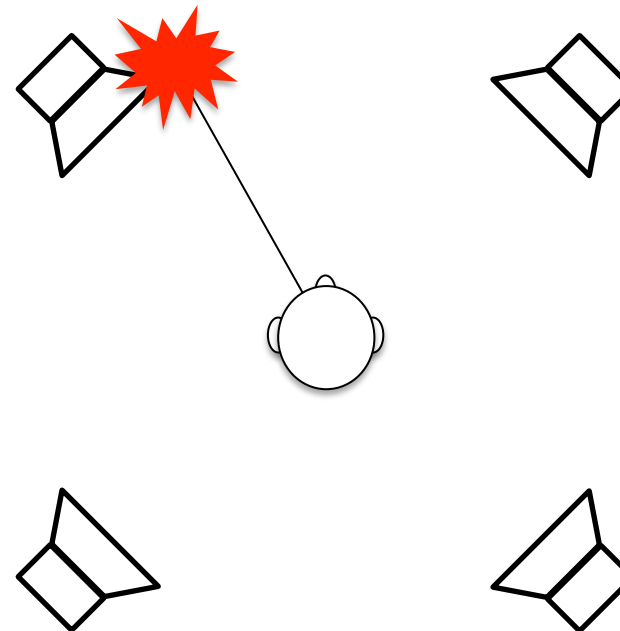
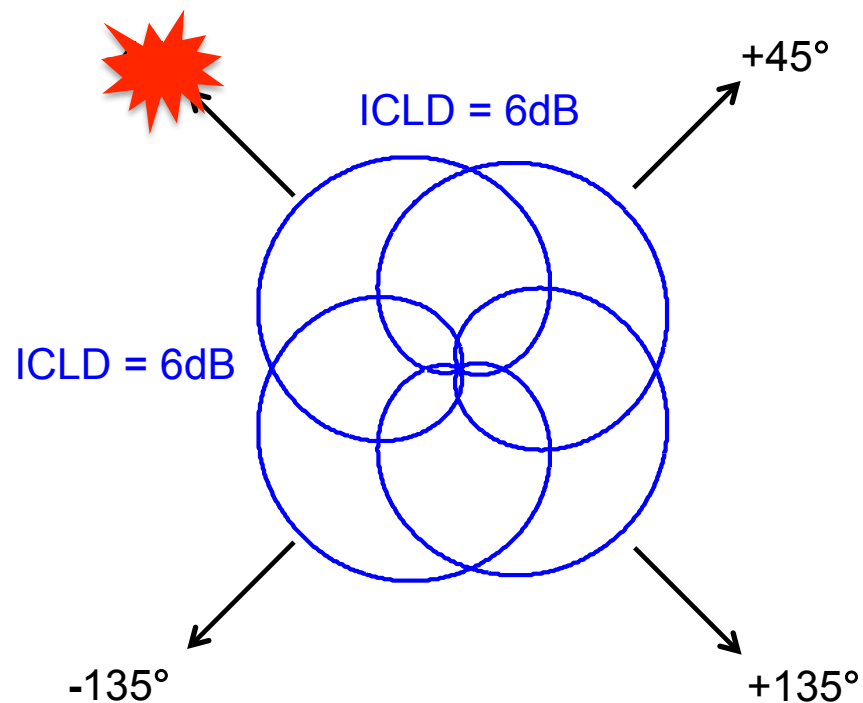
- Limitation of FOA
  - Another **6dB ICLD** for the left pair.
  - The image is perceived almost at the front left speaker (mainly one ear → no effective interaural difference)



# Psychoacoustic considerations for VR

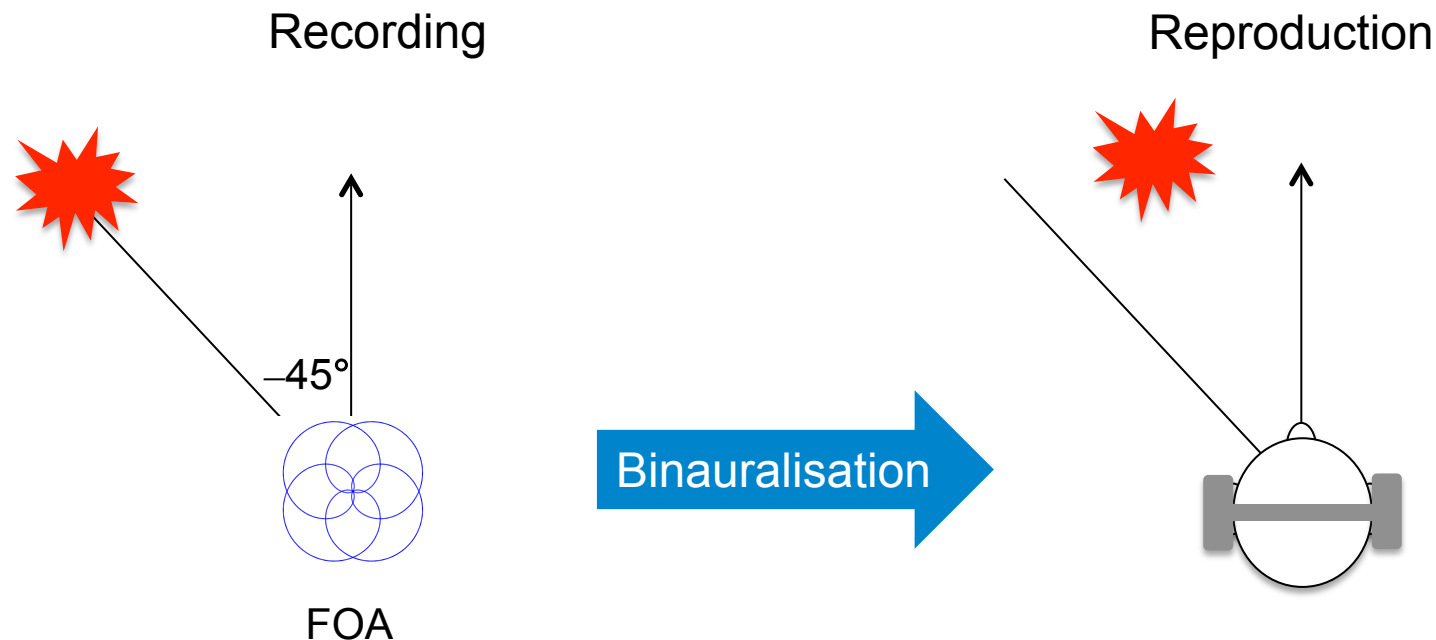
- Limitation of FOA

- The resulting image position in the quadraphonic reproduction is still not fully shifted to  $45^\circ$ .



# Psychoacoustic considerations for VR

- Problems of B-format (FOA) binauralisation for VR
  - Inaccurate localisation due to insufficient ICLD.
  - The image follows you when you rotate the head.





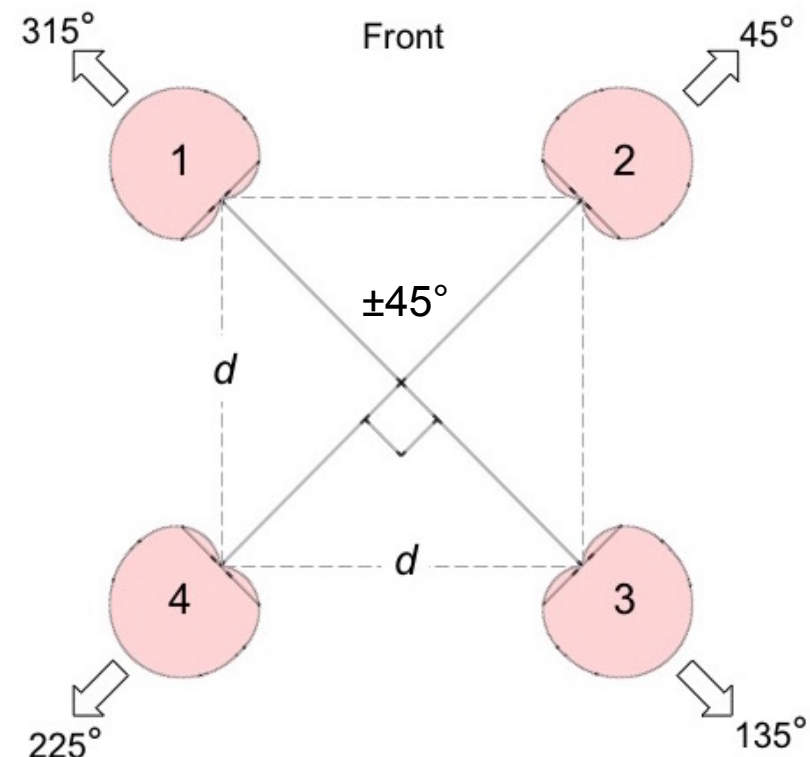
# Proposed Technique

- Equal Segment Microphone Array (ESMA)

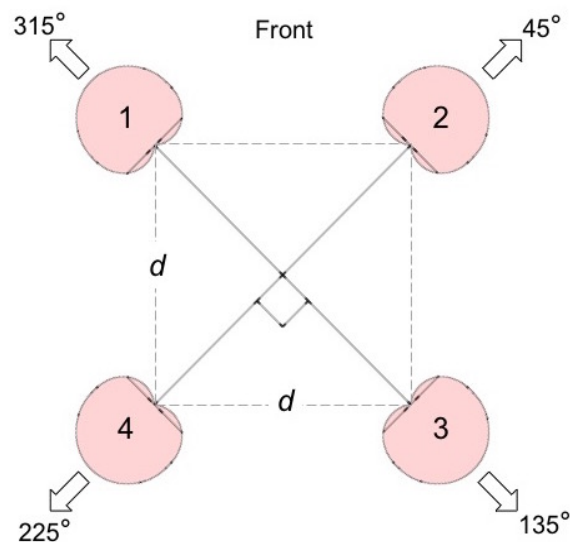
- A design concept proposed by Williams (1991), but for 360 multichannel reproduction.

- Requirements

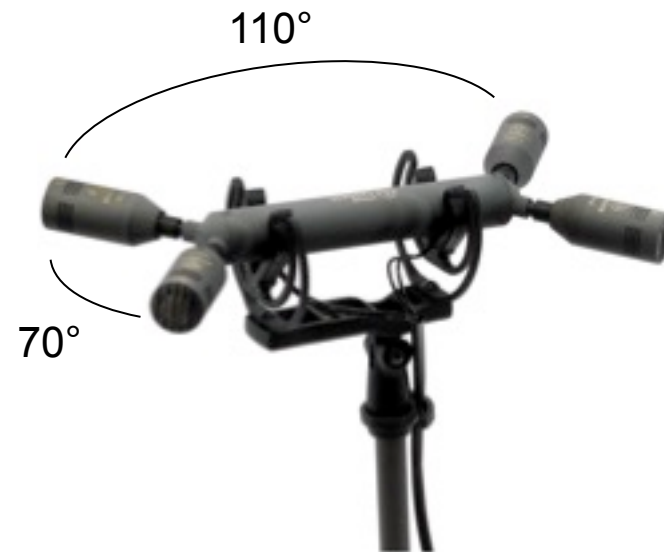
1. Equal subtended angle for all stereo segments ( $\pm 45^\circ$ ).
2. The stereophonic recording angle (SRA) of each segment should match the subtended angle of the segment. ( $\pm 45^\circ$ )



- IRT-Cross by Theile
  - Originally designed for ambience capture.
  - $d = 20$  to  $25\text{cm}$ .



- ORTF-Surround (or 3D)
  - SRA not consistent for every segment.
  - Not suitable for ESMA.

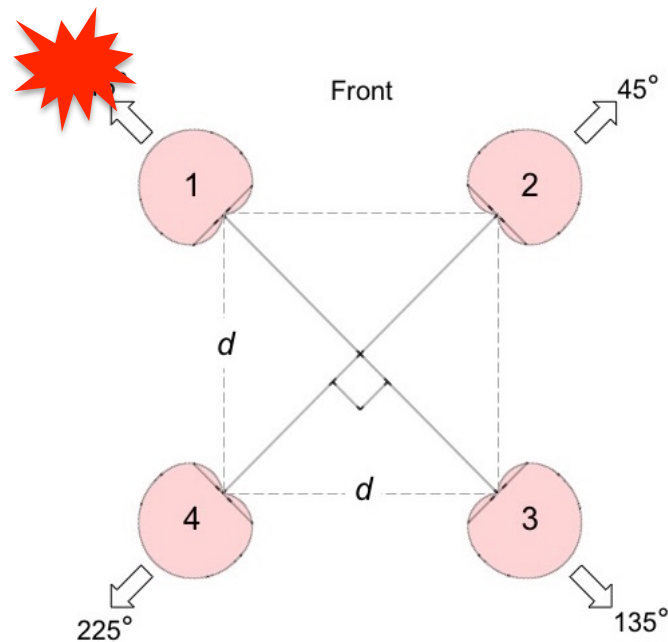


# Design philosophy

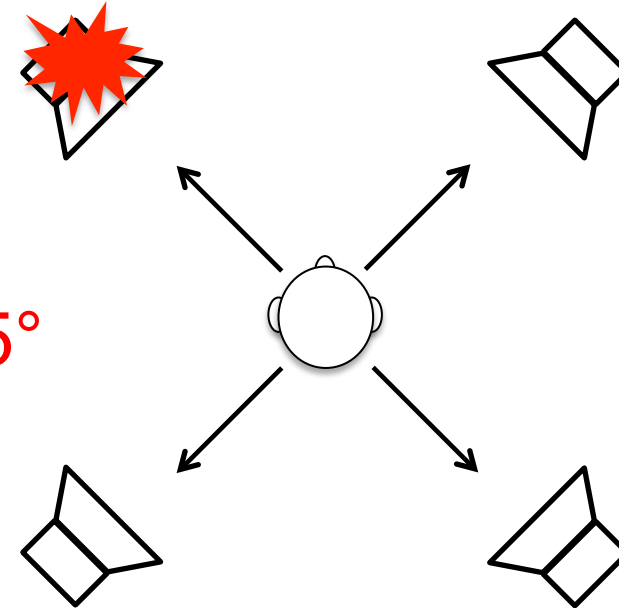
- BBC Proms using ORTF 3D



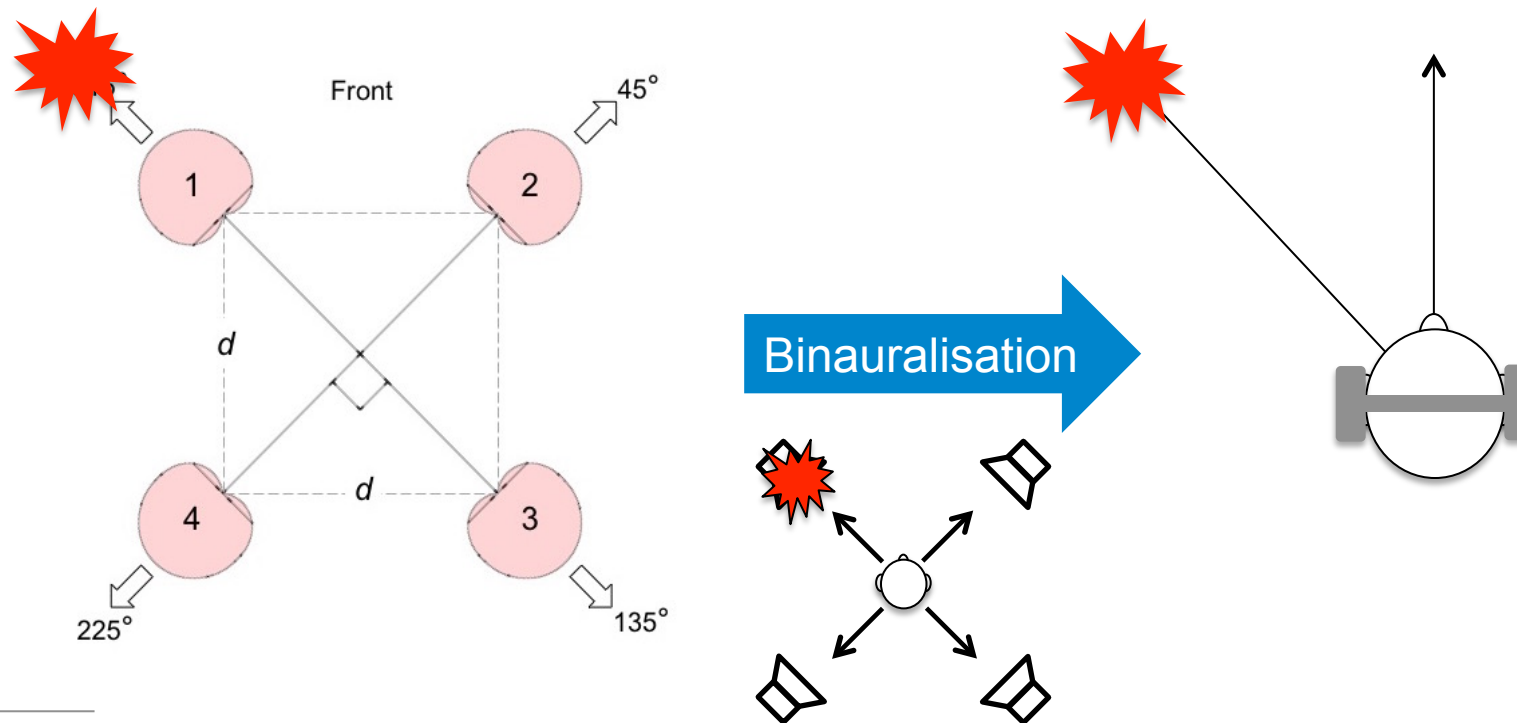
- The SRA of  $\pm 45^\circ$  for quadraphonic ESMA
  - A source at  $\pm 45^\circ$  in recording should be localised at  $\pm 45^\circ$  in reproduction.



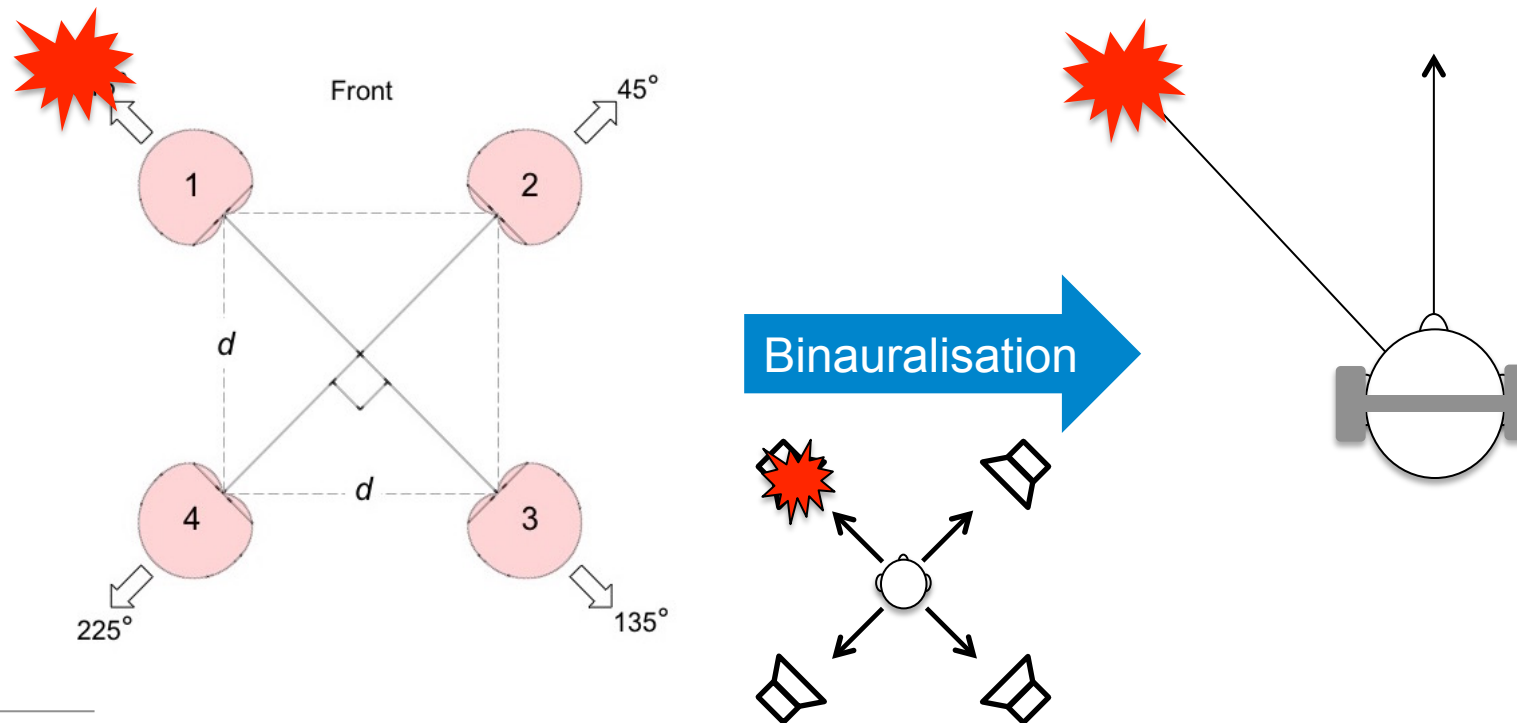
**SRA =  $\pm 45^\circ$**



- The SRA of  $\pm 45^\circ$  for quadraphonic ESMA
  - A source at  $\pm 45^\circ$  in recording should be localised at  $\pm 45^\circ$  in reproduction.



- Suitable for VR applications with head-tracking.





- The appropriate spacing between microphones to produce the  $\pm 45^\circ$  SRA?
  - Depends on what psychoacoustic time-level trade-off model is used for calculating the SRA.

Model	Microphone spacing	Source to mic array distance
Williams	23.8cm	unknown
Sengpiel	25cm	unknown
Wittek + Theile	24cm	2m
Lee + Theile	30cm	2m
Lee	50cm	2m

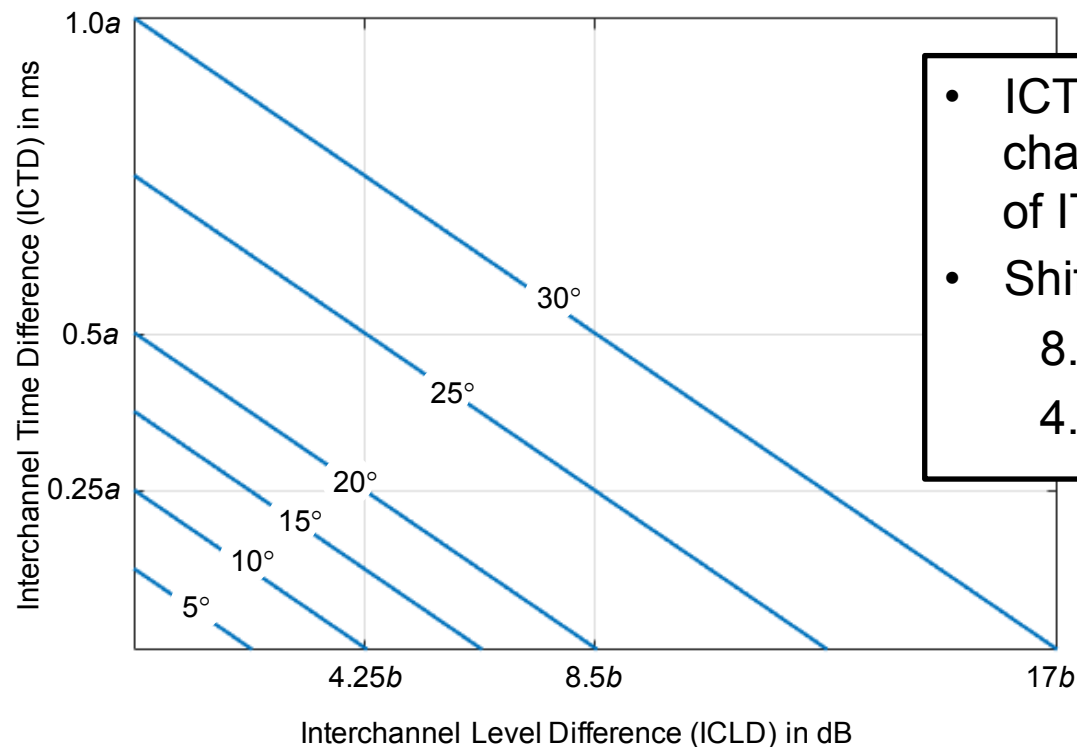
Based on ICTD and ICLD data obtained using  $\pm 30^\circ$  setup<sup>o</sup>

Optimised for  $\pm 45^\circ$  setup<sup>o</sup>



# Designing a near-coincident VR mic array

- Linear time-level trade-off functions (Lee 2016)
  - Shift region dependent.
  - Loudspeaker base angle dependent.



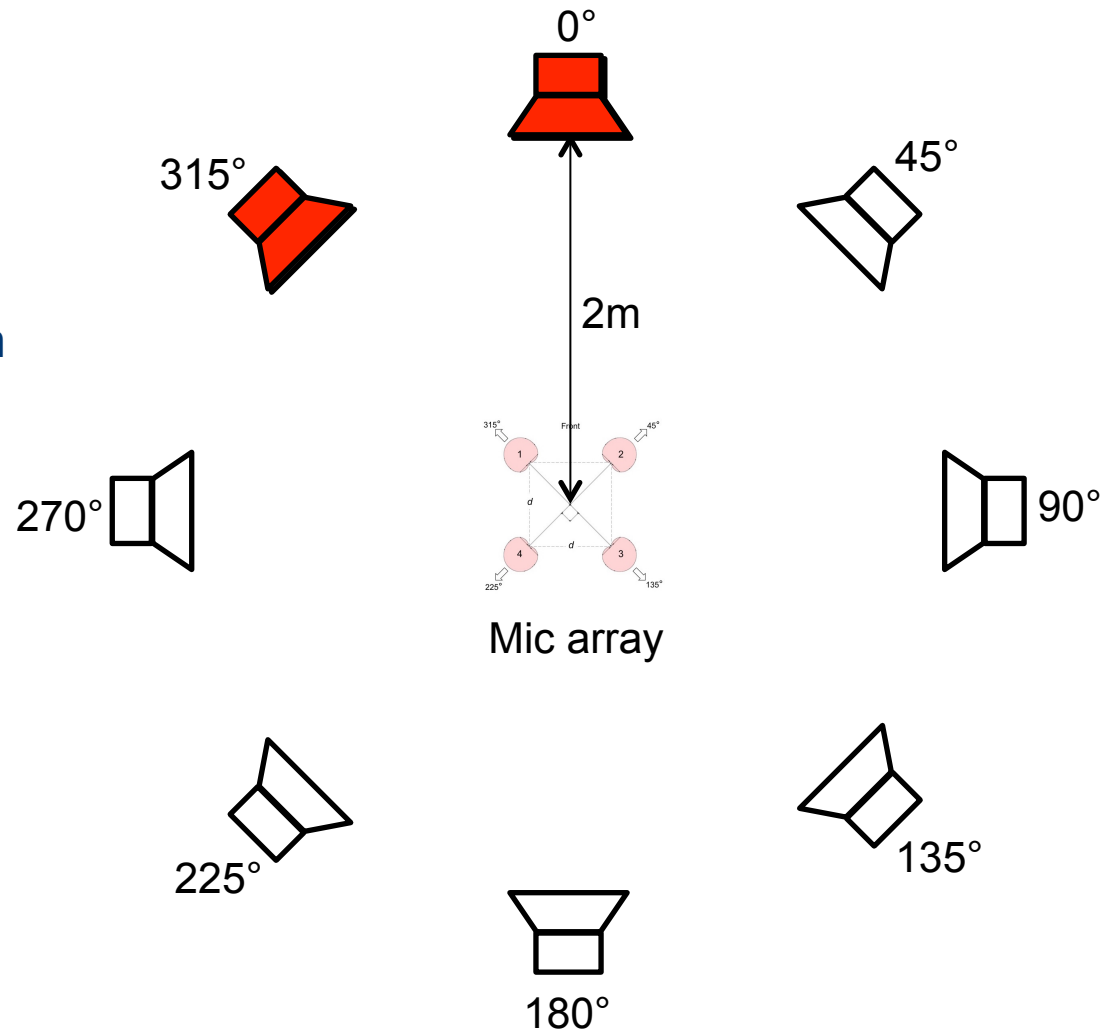
- ICTD and ICLD image shift factors change in proportion to the change of ITD and ILD.
- Shift factors for  $\pm 45^\circ$  base angle.
  - 8.8%/0.1ms; 6%/dB ( $< 30^\circ$ )
  - 4.4%/0.1ms; 3%/dB ( $30^\circ - 45^\circ$ ).

# Experiments

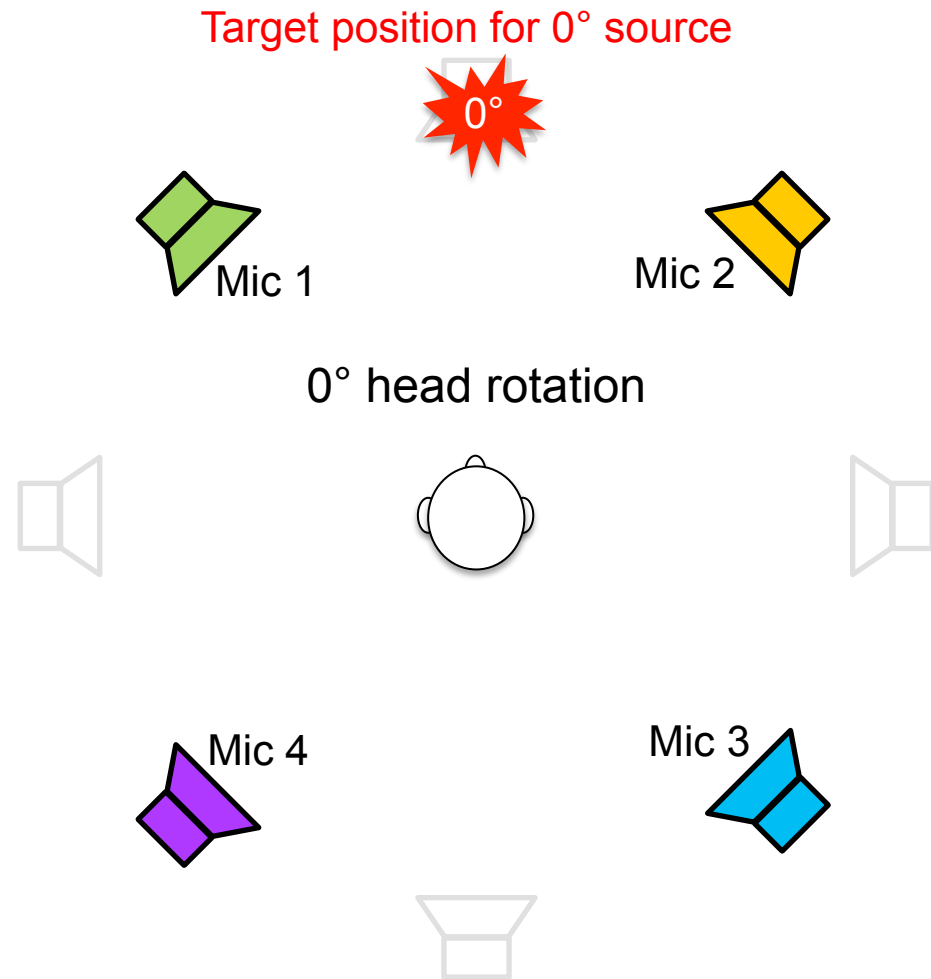
- To evaluate the localisation accuracies of the quadraphonic FOA and ESMA.
  - If the SRA of  $\pm 45^\circ$  can be achieved.
  - Loudspeaker and headphone reproduction tests in simulated head rotation scenarios.
- Microphone spacing tested:
  - 0cm (FOA)
  - 24cm (Wittek + Theile)
  - 30cm (Lee + Theile)
  - 50cm (Lee)

# Stimuli creation

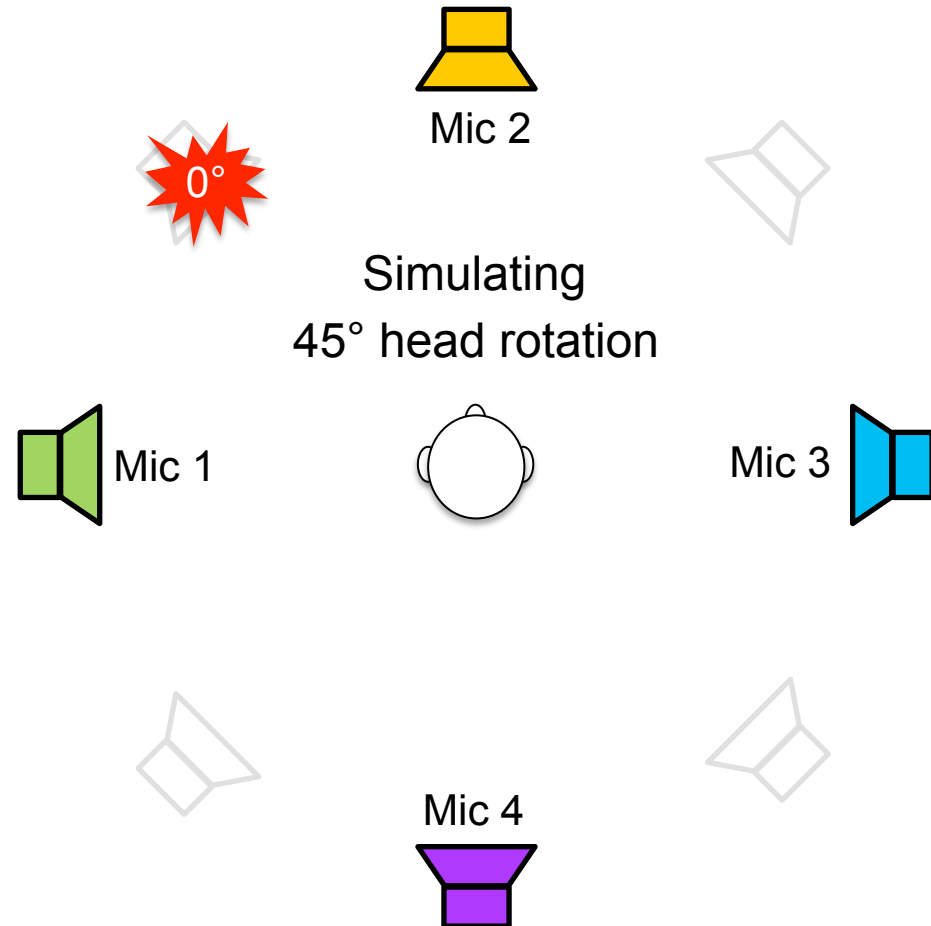
- Recording setup
  - ITU-R BS.1116 standard room.
  - 8 Genelec 8040As arranged in an octagonal layout.
  - Room impulse responses (RIRs) captured for  $0^\circ$  and  $45^\circ$ .
  - Soundfield SPS 422b for FOA.
  - Neumann KM184 for ESMA.



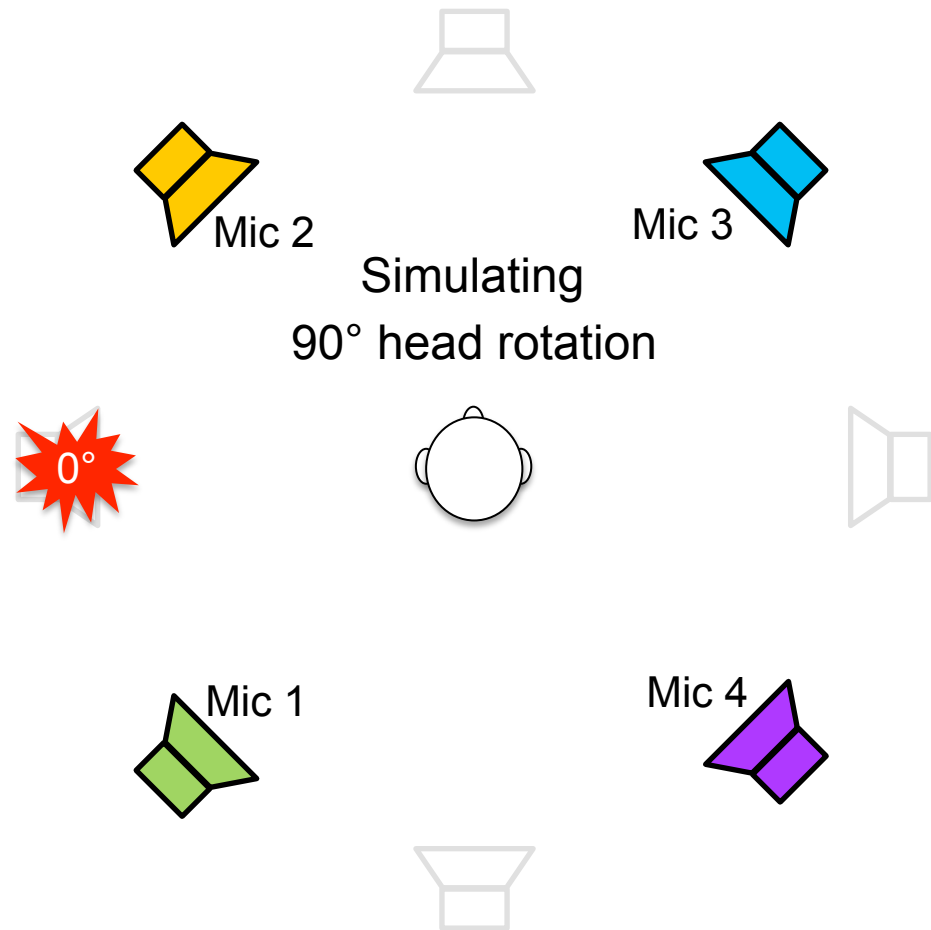
- Stimuli for Experiment 1 (Loudspeaker playback)
  - An anechoic speech signal was convolved with the direct sounds of the RIRs (reflections removed).
  - Head rotations simulated for  $0^\circ$ ,  $\pm 45^\circ$ ,  $\pm 90^\circ$ ,  $\pm 135^\circ$  and  $\pm 180^\circ$  (Soundfield rotation).



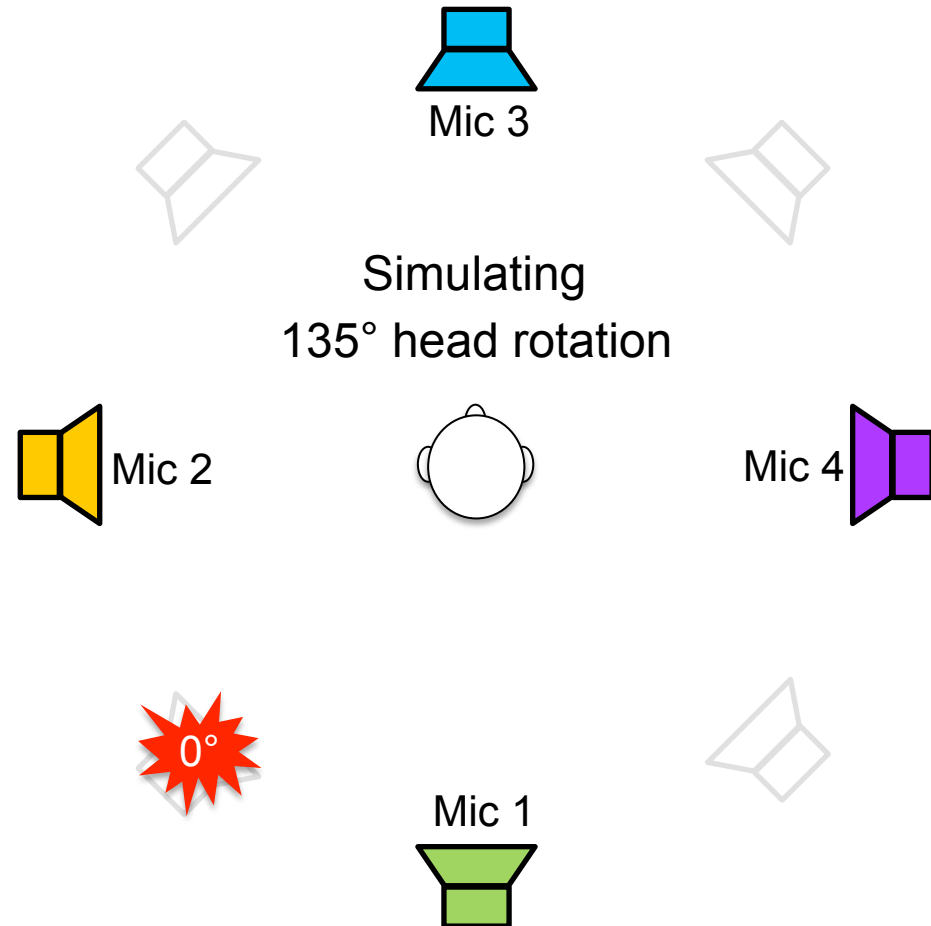
- Stimuli for Experiment 1 (Loudspeaker playback)
  - An anechoic speech signal was convolved with the direct sounds of the RIRs (reflections removed).
  - Head rotations simulated for  $0^\circ$ ,  $\pm 45^\circ$ ,  $\pm 90^\circ$ ,  $\pm 135^\circ$  and  $\pm 180^\circ$  (Soundfield rotation).



- Stimuli for Experiment 1 (Loudspeaker playback)
  - An anechoic speech signal was convolved with the direct sounds of the RIRs (reflections removed).
  - Head rotations simulated for  $0^\circ$ ,  $\pm 45^\circ$ ,  $\pm 90^\circ$ ,  $\pm 135^\circ$  and  $\pm 180^\circ$  (Soundfield rotation).

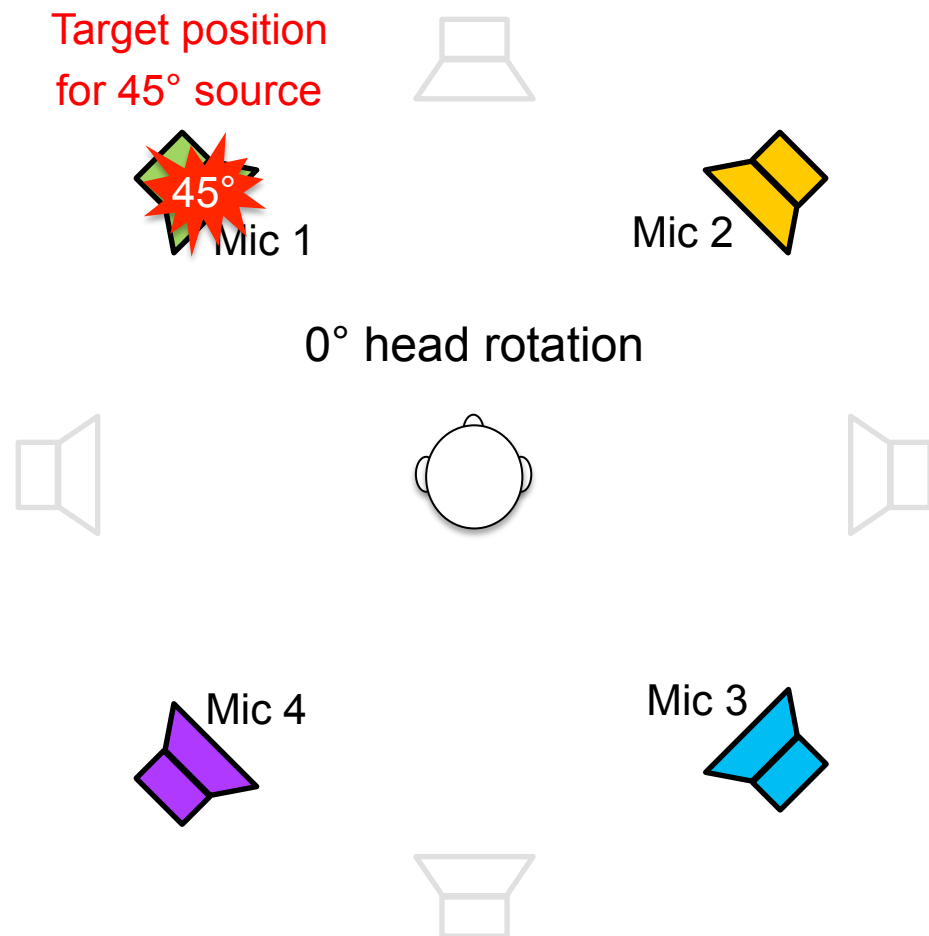


- Stimuli for Experiment 1 (Loudspeaker playback)
  - An anechoic speech signal was convolved with the direct sounds of the RIRs (reflections removed).
  - Head rotations simulated for  $0^\circ$ ,  $\pm 45^\circ$ ,  $\pm 90^\circ$ ,  $\pm 135^\circ$  and  $\pm 180^\circ$  (Soundfield rotation).

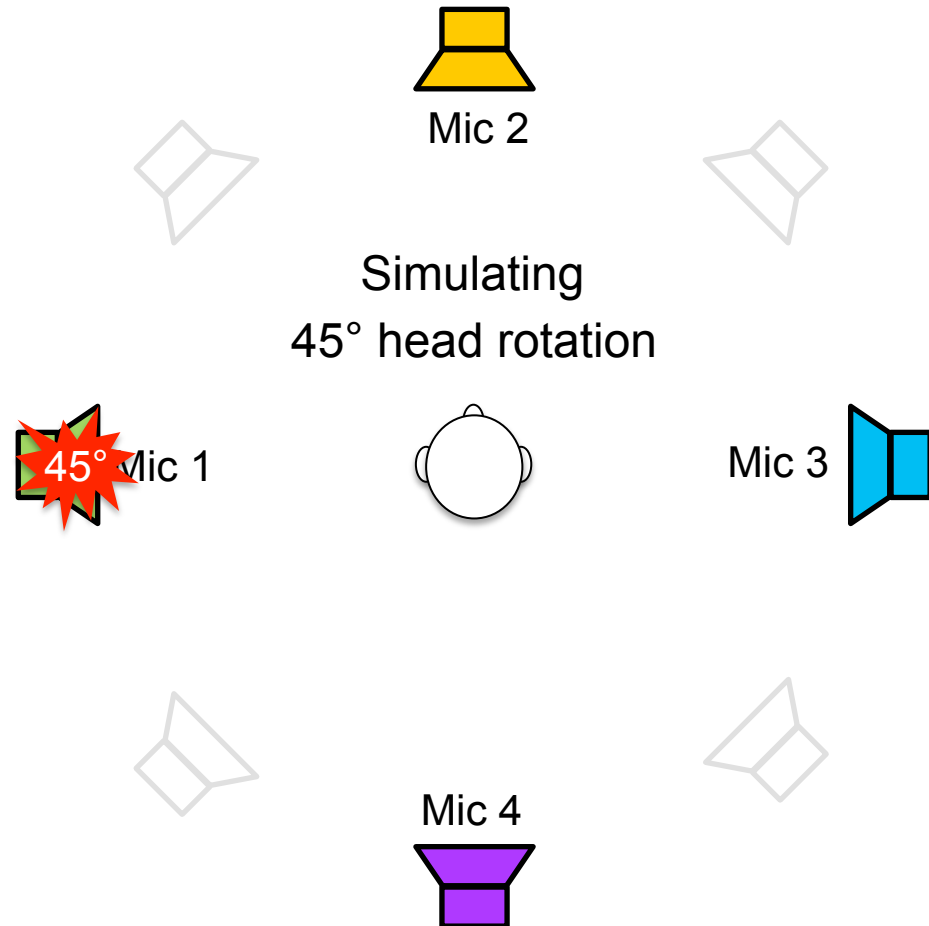




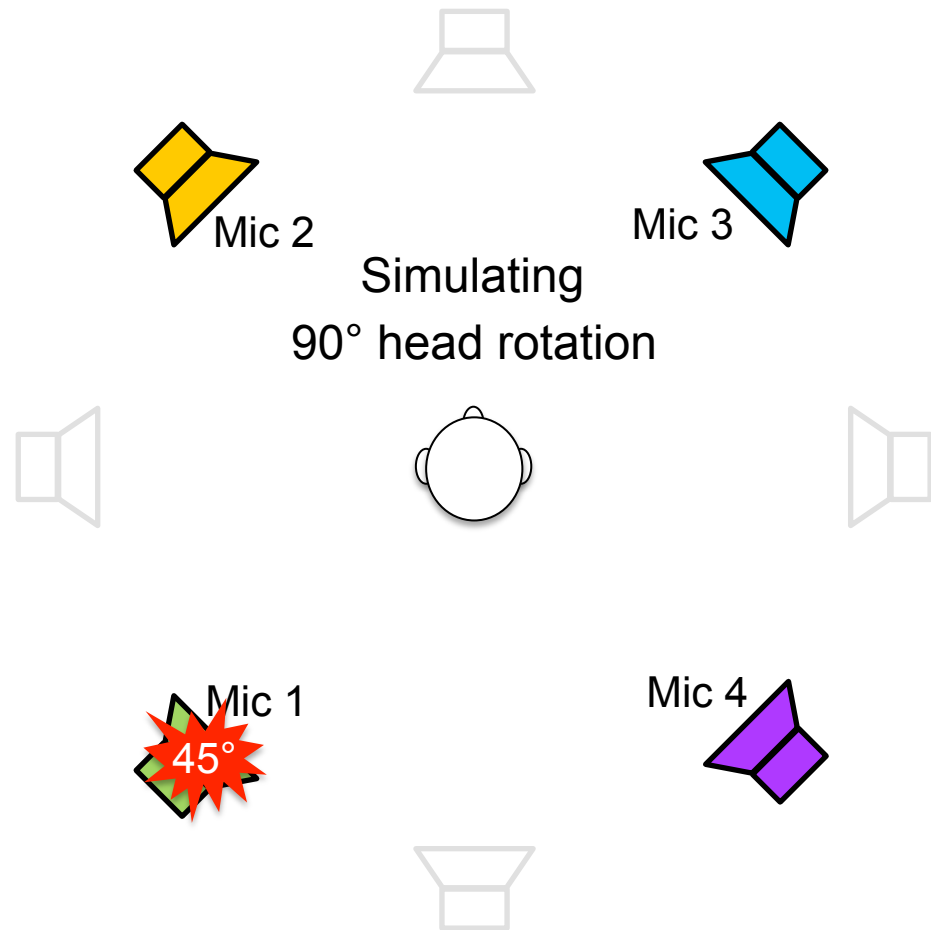
- Stimuli for Experiment 1 (Loudspeaker playback)
  - An anechoic speech signal was convolved with the direct sounds of the RIRs (reflections removed).
  - Head rotations simulated for  $0^\circ$ ,  $\pm 45^\circ$ ,  $\pm 90^\circ$ ,  $\pm 135^\circ$  and  $\pm 180^\circ$  (Soundfield rotation).



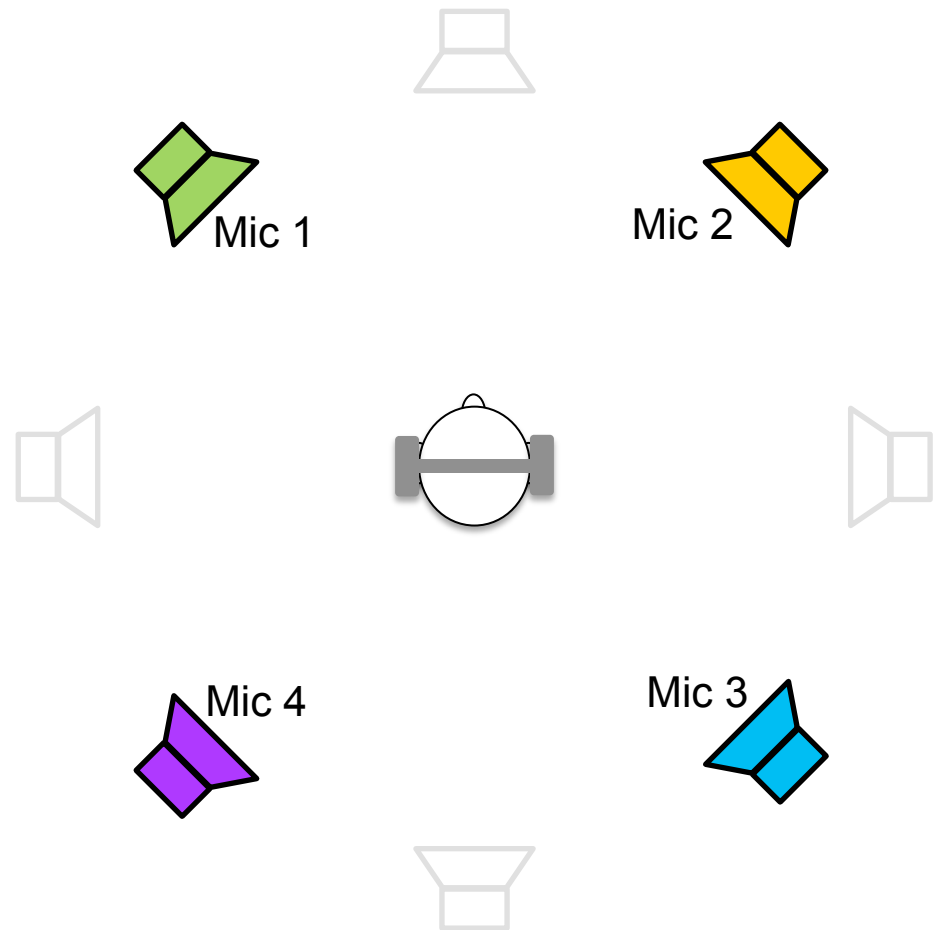
- Stimuli for Experiment 1 (Loudspeaker playback)
  - An anechoic speech signal was convolved with the direct sounds of the RIRs (reflections removed).
  - Head rotations simulated for  $0^\circ$ ,  $\pm 45^\circ$ ,  $\pm 90^\circ$ ,  $\pm 135^\circ$  and  $\pm 180^\circ$  (Soundfield rotation).



- Stimuli for Experiment 1 (Loudspeaker playback)
  - An anechoic speech signal was convolved with the direct sounds of the RIRs (reflections removed).
  - Head rotations simulated for  $0^\circ$ ,  $\pm 45^\circ$ ,  $\pm 90^\circ$ ,  $\pm 135^\circ$  and  $\pm 180^\circ$  (Soundfield rotation).

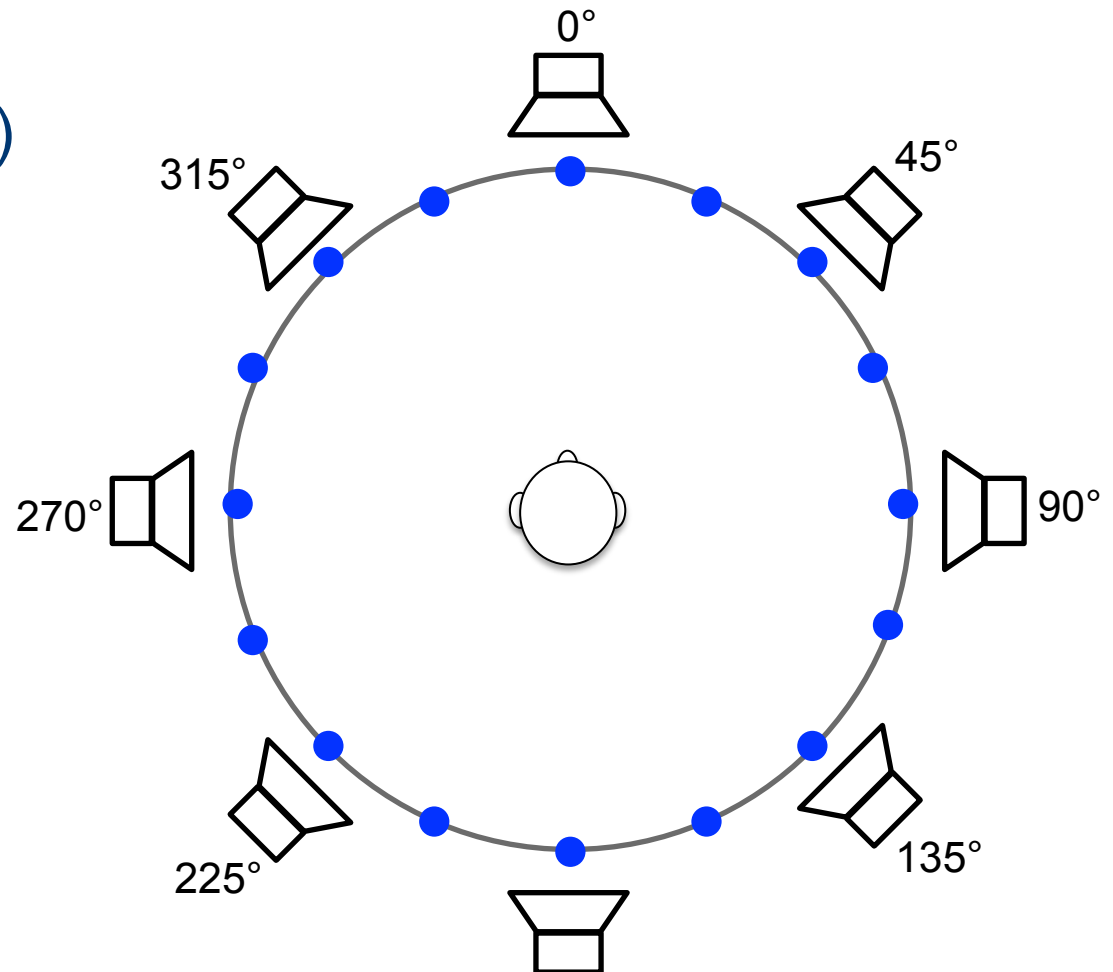


- Stimuli for Experiment 2  
(Binaural playback)
  - Same conditions as Experiment 1, but with the full RIRs (reflections included).
  - The multichannel stimuli were binauralised with dry KU100 dummy head HRIRs from the ‘SADIE’ database (Kearney 2015).



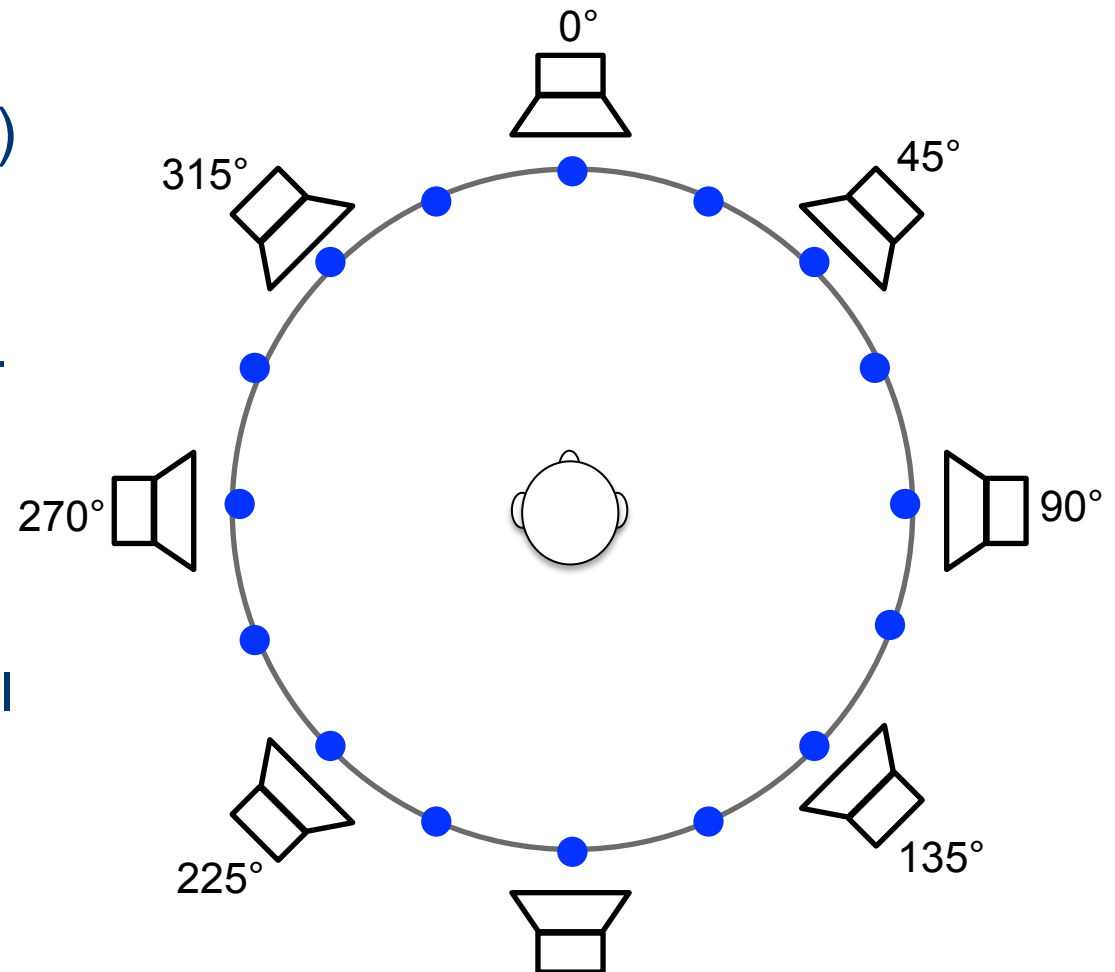
# Listening tests

- Experiment 1  
(Loudspeaker playback)
  - Loudspeakers hidden by acoustically transparent curtains.
  - Small markers were placed on the curtain from  $0^\circ$  with  $22.5^\circ$  intervals.
  - 70dBA playback level.



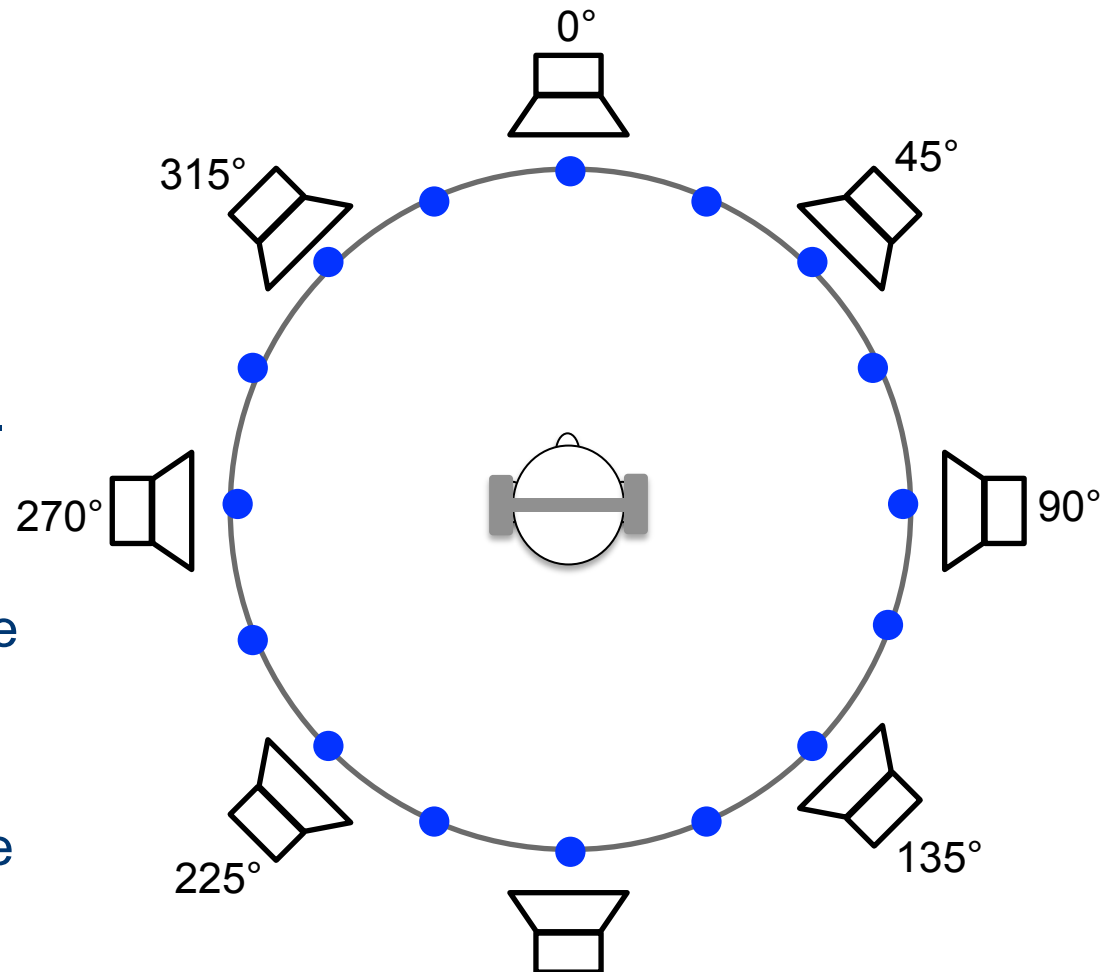
# Listening tests

- Experiment 1  
(Loudspeaker playback)
  - 9 experienced subjects repeated each test twice.
  - The task was to mark down the perceived image position on a horizontal circle on a GUI with markers indicated with  $22.5^\circ$  intervals.

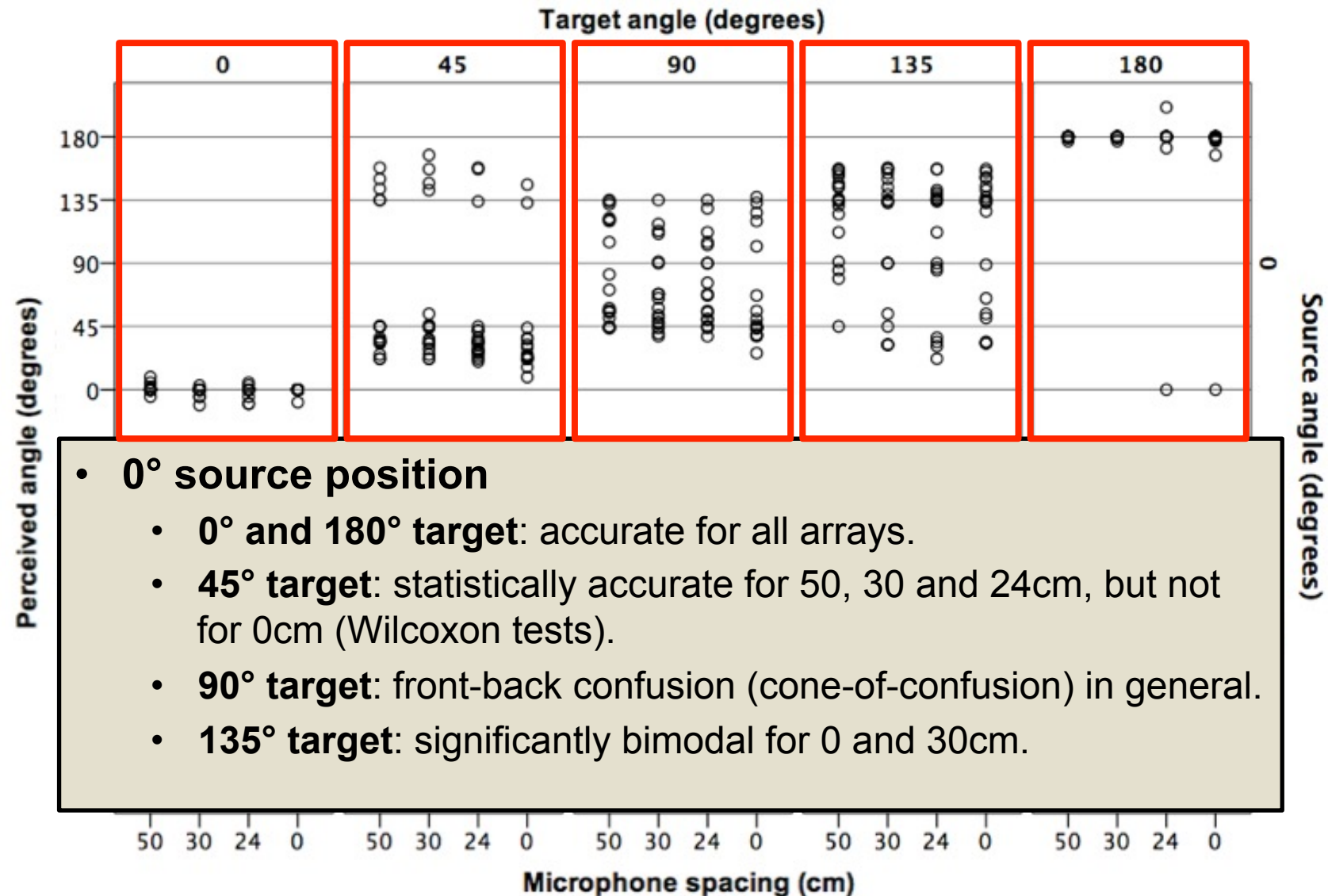


# Listening tests

- Experiment 2  
(Binaural playback)
  - The same room, subjects, task and method as Experiment 1.
  - Equalised Sennheiser HD650 headphones were used.
  - Loudness matched to the playback levels of multichannel stimuli.

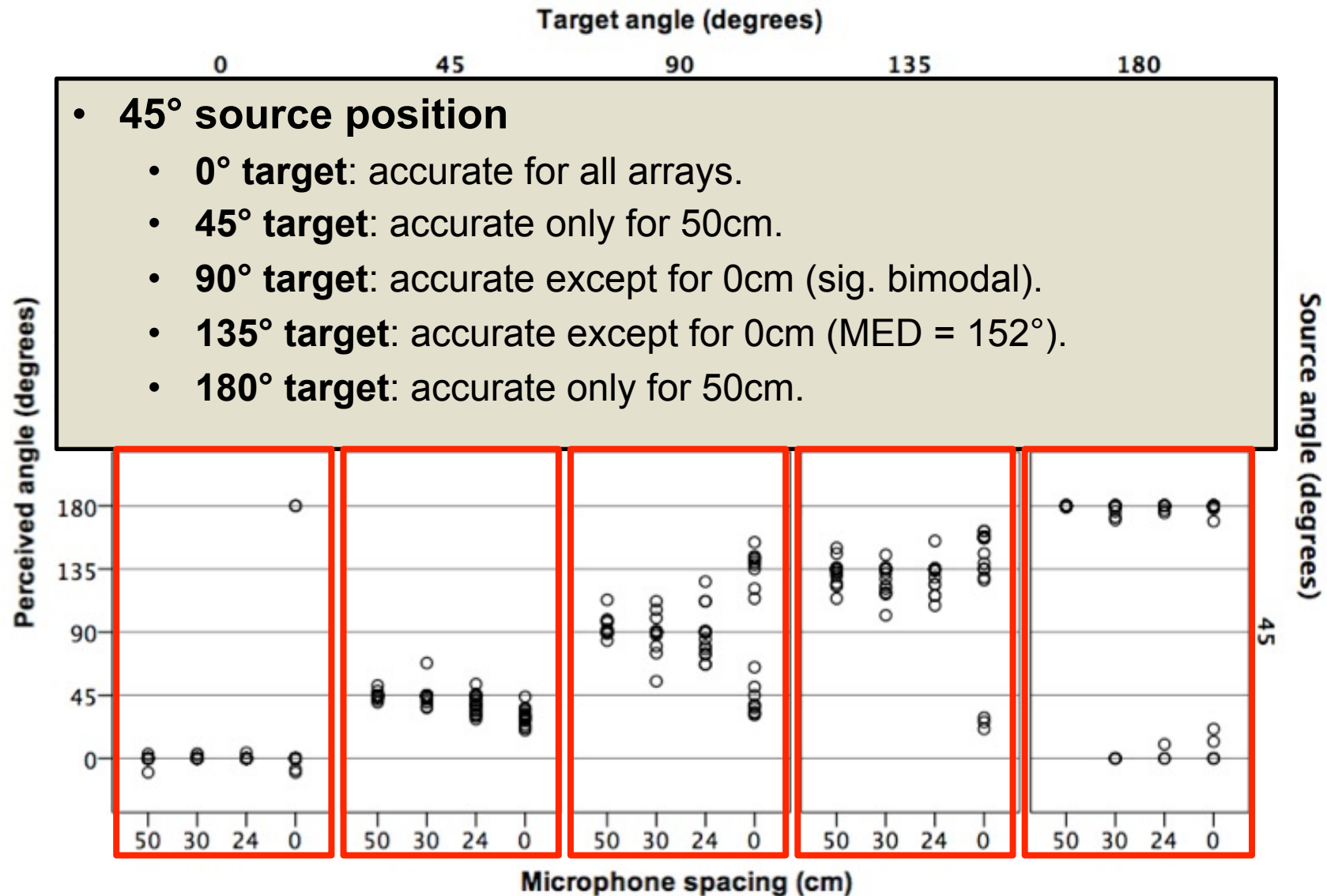


# Results – Loudspeaker experiment

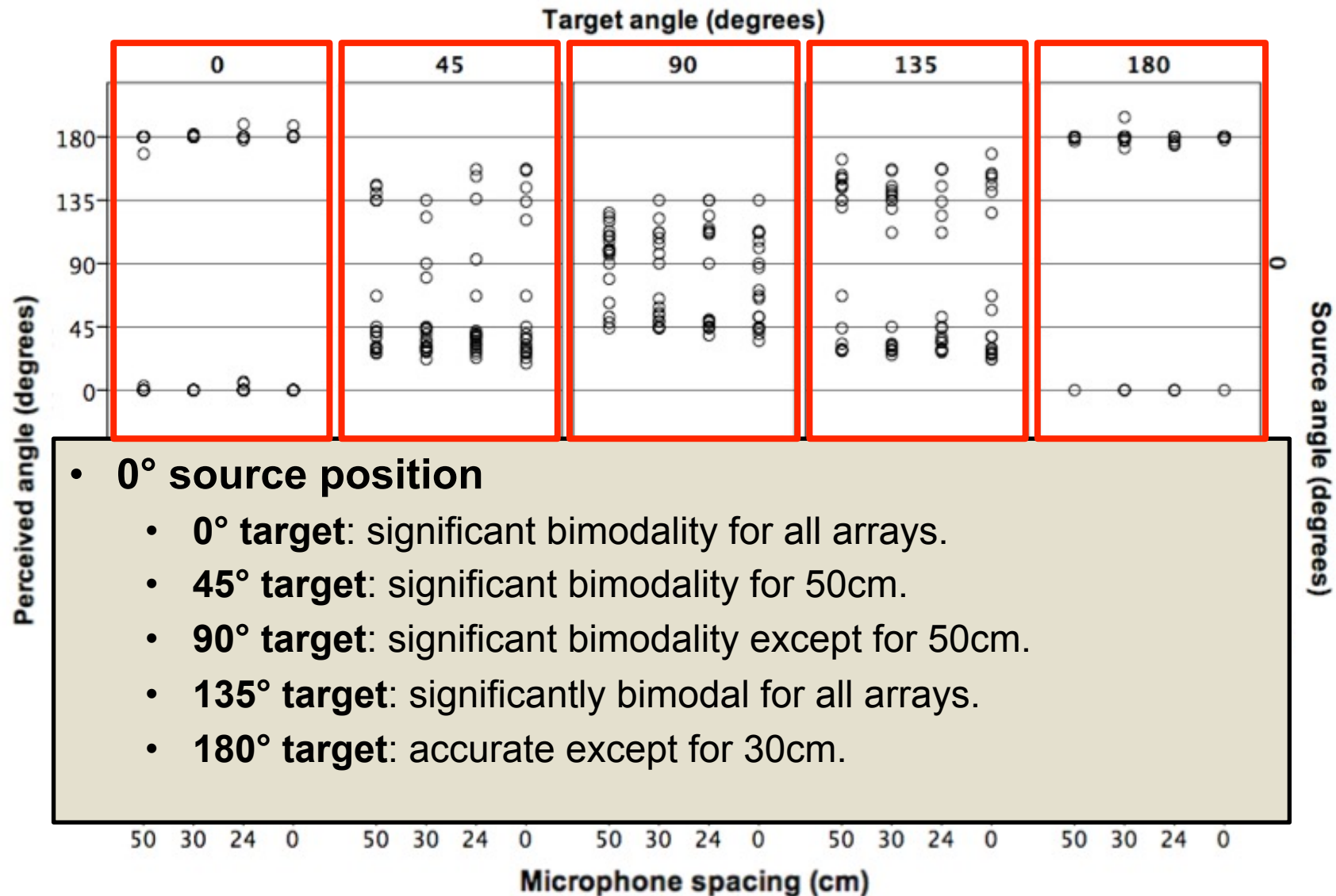




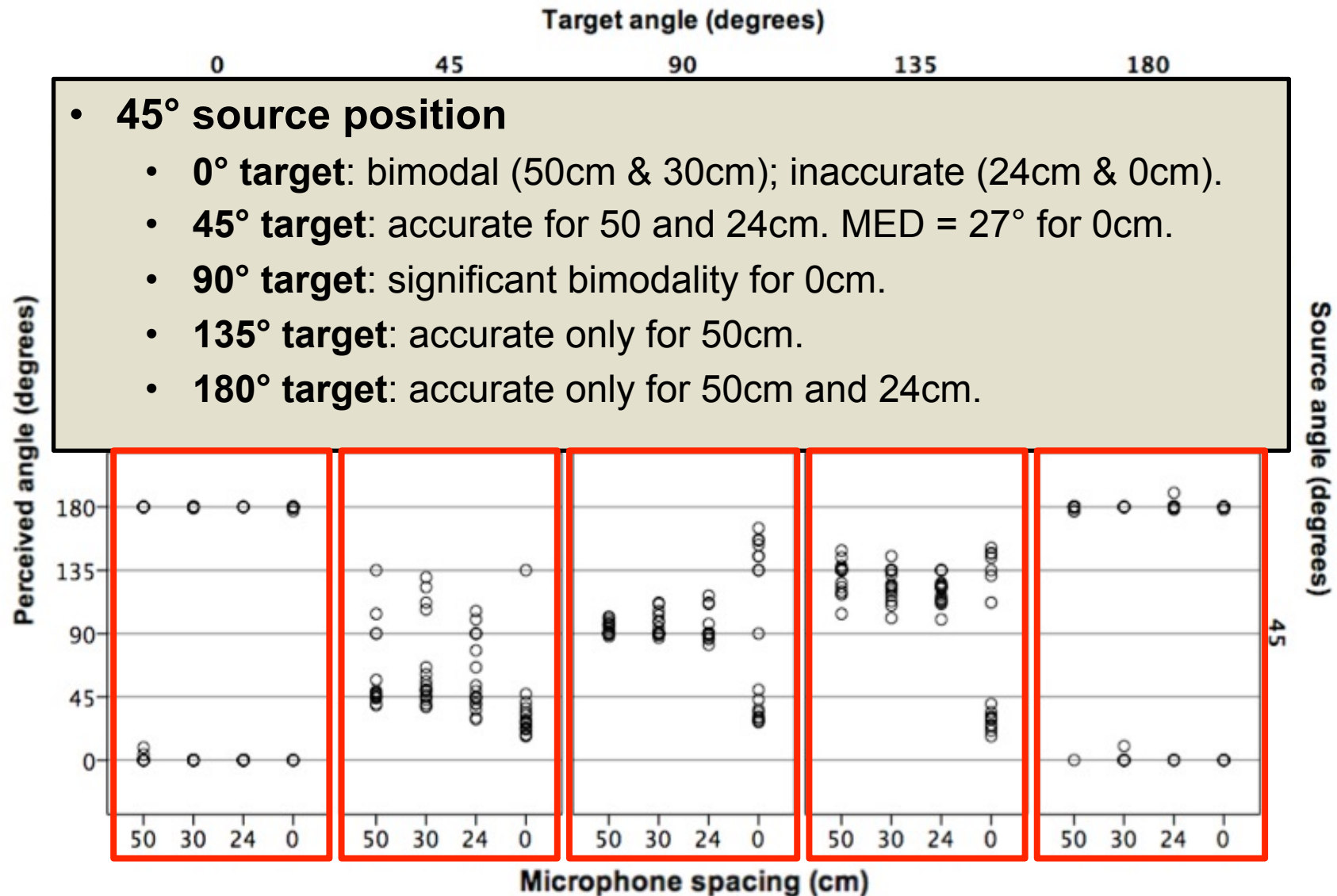
# Results – Loudspeaker experiment



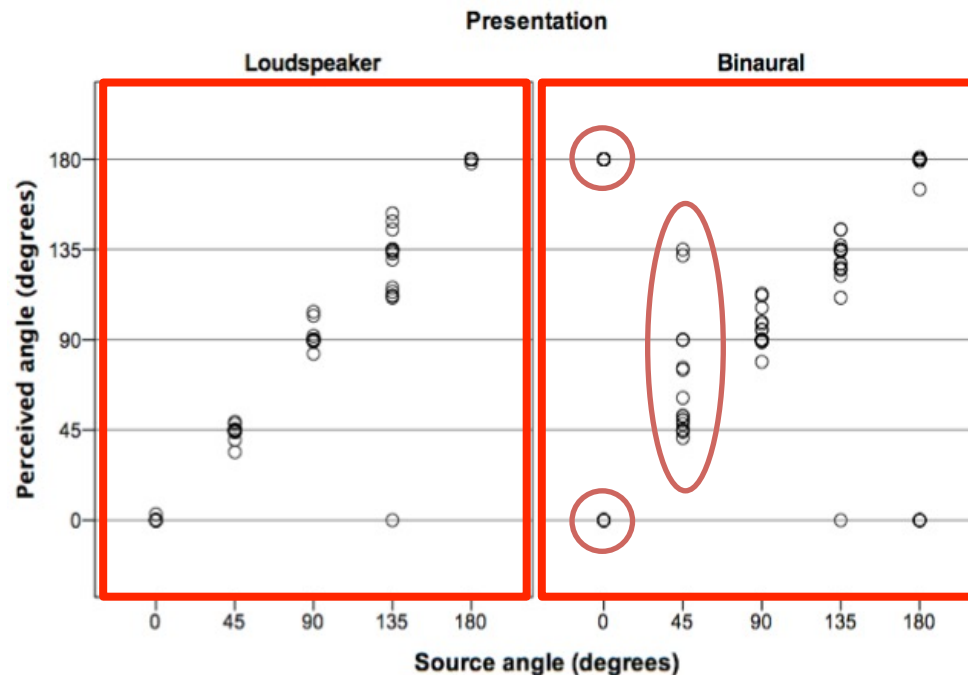
# Results – Binaural experiment



# Results – Binaural experiment



# Results – Real source



- **Loudspeaker**

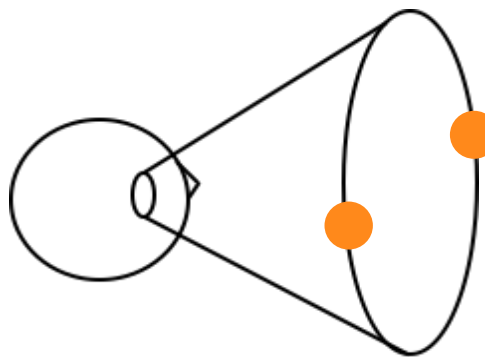
- Loudspeaker: accurate for all source angles.

- **Binaural**

- Binaural responses are generally more spread than loudspeaker ones.
- 0°: significantly bimodal.
- 45°: inaccurate, MED = 52°.
- 90°, 135°: accurate.
- 180°: inaccurate, bimodal.

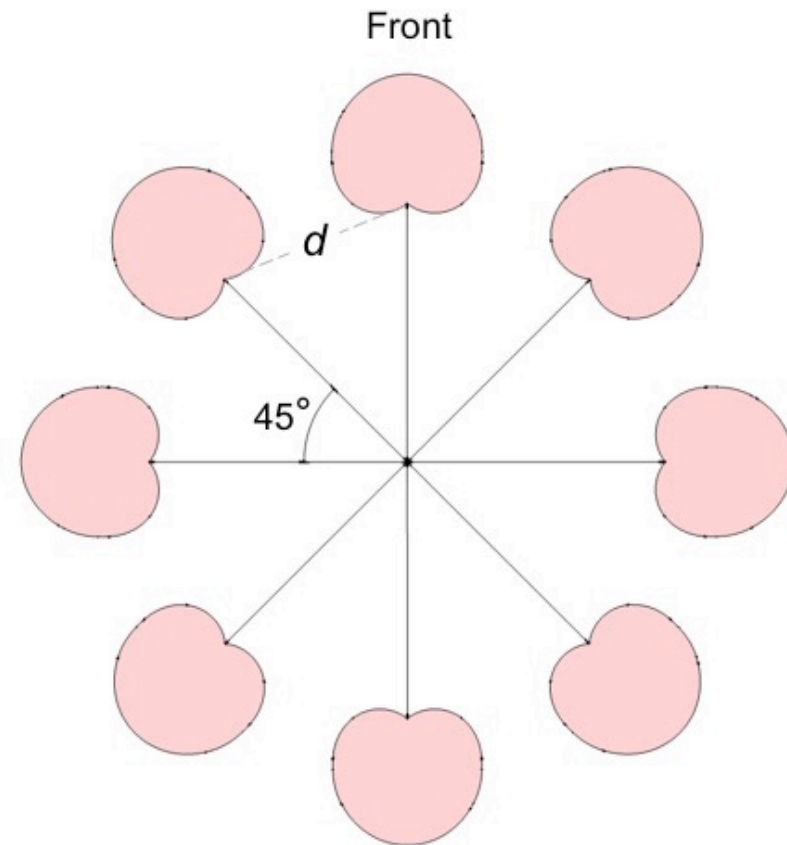
- Microphone spacing effect
  - 0cm had the worst localisation performance overall.
    - Significant bimodal distributions for many target angle conditions.
    - Perceived to be significantly narrower for the 45° source in both loudspeaker (MED = 30°) and binaural (MED = 27°).
  - 50cm was the only spacing that achieved the SRA of  $\pm 45^\circ$ .
    - Seems to validate the new psychoacoustic model.
  - 50cm had slightly better consistency and accuracy than the other configurations overall.
    - But a smaller size might be more beneficial in practical situations.
    - Practical importance of localisation accuracy in VR?

- Source angle effect
  - The  $0^\circ$  source produced larger response spreads and more bimodal distributions than the  $45^\circ$ .
    - Front-back confusion (Cone of confusion), especially for the  $90^\circ$  target angle.
    - Lateral phantom image localisation is highly unstable (Theile and Plenge 1977, Martin et al 1999).



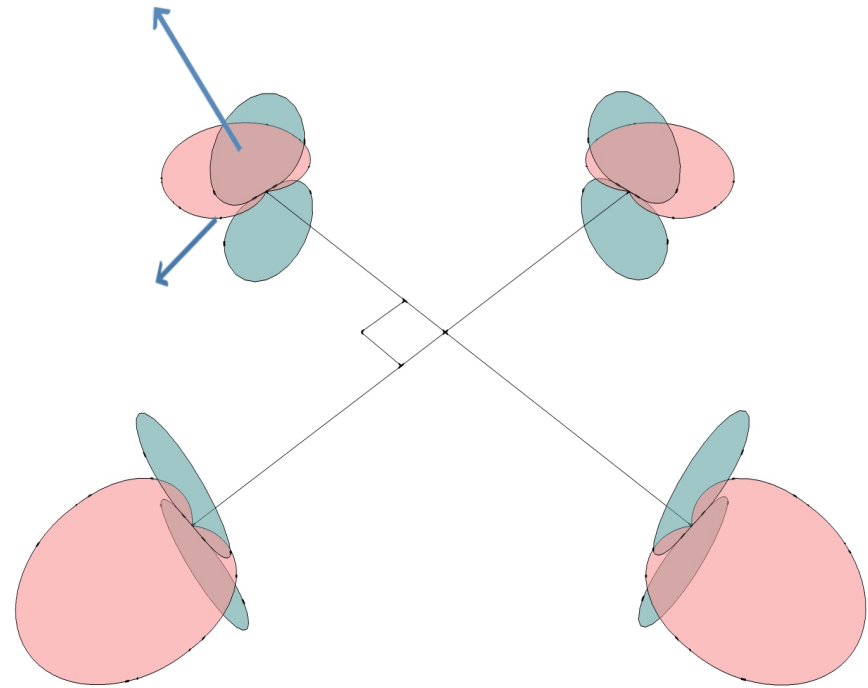
- Loudspeaker vs. Binaural
  - **Front-back confusion** was more frequently observed in the binaural presentation, but not in the loudspeaker one.
  - The binaural presentation had more spread responses.
  - Real source results also show similar tendencies for the 0° and 45°.
  - Might be due to the use of **non-individualised HRTF**, rather than the microphone arrays.
  - But more about the **lack of head movement?**
    - FB confusion can occur even with individualised HRTF when head rotation is not allowed (Wightman and Kistler 1999).
  - The FB confusion problem might be largely resolved in practical VR applications with head tracking.

- Higher Order ESMA
  - For an octagonal setup, each segment should have the SRA of  $45^\circ$  ( $\pm 22.5^\circ$ ).
  - Can potentially solve the problem of unstable side image localisation.
  - Mic spacing  $d$ 
    - *Williams: 82cm*
    - *Lee: 55cm*





- Adding vertical dimension to ESMA
  - Cardioid + Figure-of-eight in a vertically coincident fashion.
    - Vertical Mid-Side decoding.
    - Vertical microphone spacing has little effect on LEV (Lee and Gribben JAES 2014).
    - Vertical level panning can provide source imaging with a limited resolution (Barbour 2003, Mironovs and Lee 2016).
    - Vertical time panning is highly unstable (Wallis and Lee JAES 2015).



- ESMAAs had a better localisation accuracy than FOA.
- 50cm spacing had the best localisation accuracy, but 30cm or 24cm might still be acceptable.
- Front-Back confusion in binaural reproduction without head rotation.
- Ongoing works
  - Investigations on different attributes.
  - Externalisation, tonal quality, spaciousness, naturalness, etc.
  - Practical evaluations with head tracking.

Thank you for listening.

Hyunkook Lee  
[h.lee@hud.ac.uk](mailto:h.lee@hud.ac.uk)